# EEC 289A Assignment 1 Report
**Chenye Yang, Hanchu Zhou, Haodong Liang, Yibo Ma**

## 1 Introduction

In this project, we aim to do K-mean clustering on the patches from MNIST dataset, shown in Figure 1. After the clustering, we observe the results, and try to answer the following interesting questions:

1. What is the change of the learned clusters when K increases from 100 - 10,000?

2. How well one can reconstruct a 5x5 MNIST patch by the learned dictionary (clusters)?

3. How many clusters does one need in order to cover the whole patch space?

4. What are these clusters and do they have any interpretable meanings?

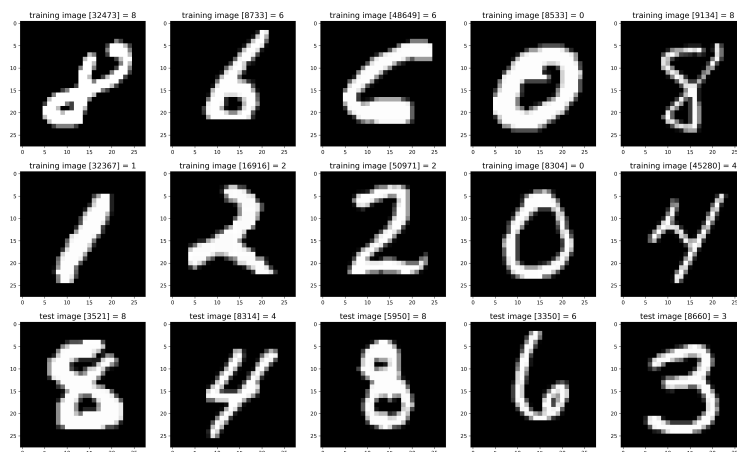5. How is one digit made from these clusters?



Figure 1: Some examples of the MNIST dataset

## 2 Methodology

### 2.1 Data Preprocessing

We load the MNIST dataset ($28 \times 28$ handwritten digits)[1], including 60000 training images and 10000 testing images. The pixel values in each image are not normalized, and they range from 0 (black) to 255 (white). Then, we extract the patches from the training images, by sliding a $5 \times 5$ window over the images. Therefore, for each handwritten digit, we will get $(28 - 5 + 1) \times (28 - 5 + 1)$ patches. Then we get rid of all the blank patches, leading to 20,074,704 non-blank patches in total. Each patch is reshaped to a 25-dimensional vector, and we get a matrix $X \in R^{20,074,704 \times 25}$.

### 2.2 K-mean clustering

Once we get all the patches, we can do the K-mean clustering on the patches. The K-mean clustering is a method to partition the data into K clusters, where each data point belongs to the cluster with the nearest mean. Although there is a wildly used K-means algorithm in library *scikit learn*[2], we also implement the

---

[1]Downloaded from https://git-disl.github.io/GTDLBench/datasets/mnist_datasets/
[2]https://scikit-learn.org/stable/modules/generated/sklearn.cluster.KMeans.html

K-means clustering algorithm using PyTorch, hoping to achieve acceleration by using Nvidia CUDA[3] or Apple MPS[4], which is shown in the Algorithm 1.

We define that:

- $K$: Number of the clustering

- $C$: Centroids of the clustering

- $P$: The $5 \times 5$ patch

- $X$: The entire dataset of patches with size $20,074,704 \times 25$

---

**Algorithm 1** PyTorch K-means Clustering for $5 \times 5$ Patches

---

1: **procedure** KMEANSPATCHES($Patches, K$)
2:     $X \leftarrow$ Reshape each patch $P$ in $Patches$ to $R^{25}$                    ▷ No normalization
3:     Initialize $C$ random $K$ data points from $X$
4:     $C_{\text{old}} \leftarrow$ Copy of $C$
5:     **repeat**
6:         Compute distances from each vector $X_i$ in $X$ to each centroid in $C$
7:         Assign each vector $X_i$ to the closest centroid
8:         **for** $j \leftarrow 1$ to $K$ **do**
9:             **if** Count($X_i$ assigned to $C_j$) $= 0$ **then**
10:                 Randomly reinitialize centroid $C_j$ from $X$
11:             **else**
12:                 Update centroid $C_j$ by calculating the mean of all vectors assigned to $C_j$
13:             **end if**
14:         **end for**
15:         $C_{\text{new}} \leftarrow$ Copy of $C$
16:         $C_{\text{move}} \leftarrow \text{norm}(C_{\text{new}} - C_{\text{old}})$
17:         $C_{\text{old}} \leftarrow C_{\text{new}}$
18:     **until** $C_{\text{move}} <$ tolerance
19:     **return** Updated centroids and cluster labels
20: **end procedure**

---

## 2.3   Reconstruction

After the clustering, we can use the centroids to reconstruct the original digits, as shown in Algorithm 2. The reconstruction is done by assigning each non-blank patch to the nearest centroid, and keep the blank patches as zeros. Because of the overlap of the patches, we will have multiple centroids assigned to the same pixel. Thus, we need a count matrix to record the number of how many centroids are assigned to each pixel. Finally, we average the value of each pixel to get the reconstructed digit.

Some notations in this algorithm:

- $POS$: The position indices of a specific patch

- $D$: The ground truth of the handwritten digit

- $\hat{D}$: The reconstructed result of $D$ according to the K-mean clustering centroids

---

[3]`https://pytorch.org/docs/stable/cuda.html`
[4]`https://developer.apple.com/metal/pytorch/`

---

**Algorithm 2** Reconstruct Handwritten Digit Images

---

1: **procedure** RECONSTRUCTDIGIT($D, K, Model$)
2:     $P \leftarrow$ empty list to store patches
3:     $Pos \leftarrow$ empty list to store position indices
4:     **for** $i \leftarrow 1$ to 24 **do**                                                                        ▷ 28 - 5 + 1 = 24
5:         **for** $j \leftarrow 1$ to 24 **do**
6:             $patch \leftarrow D[i : i + 5, j : j + 5]$
7:             **if** not all zeros in $patch$ **then**
8:                 $P$.append($patch.flatten()$)
9:                 $Pos$.append($(i, j)$)
10:             **end if**
11:         **end for**
12:     **end for**
13:     $X \leftarrow$ stack of patches in $P$
14:     $Labels \leftarrow Model$.predict($X$)                                                ▷ Assign patches to centroids
15:     $\hat{D} \leftarrow$ zero matrix of the same size as $D$                       ▷ Initialize reconstructed digit with zeros
16:     $Count \leftarrow$ zero matrix of the same size as $D$                          ▷ Initialize Count matrix with zeros
17:     **for** $k \leftarrow 0$ to len($P$) $- 1$ **do**
18:         $pos \leftarrow Pos[k]$
19:         $cluster \leftarrow Labels[k]$
20:         $centroid \leftarrow Model.cluster\_centers_[cluster].reshape(5, 5)$
21:         $\hat{D}[pos[0] : pos[0] + 5, pos[1] : pos[1] + 5] \leftarrow \hat{D}[pos[0] : pos[0] + 5, pos[1] : pos[1] + 5] + centroid$
22:         $Count[pos[0] : pos[0] + 5, pos[1] : pos[1] + 5] \leftarrow Count[pos[0] : pos[0] + 5, pos[1] : pos[1] + 5] + 1$
23:     **end for**
24:     $Count[Count == 0] = 1$                                                                ▷ Avoid division by 0
25:     **return** $\hat{D}/Count$
26: **end procedure**

---

# 3 Experiment Results

## 3.1 Clustering Results with Different $K$

We conducted the K-mean clustering on the patches with different $K$ values, where:

$$K = 100, 200, 300, \ldots, 900, 1000, 2000, \ldots, 9000, 10000.$$

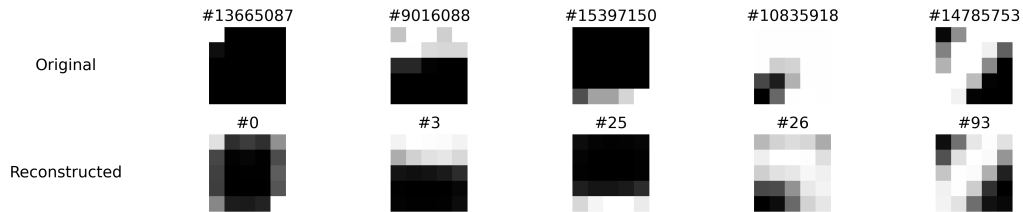Some results of the patches and the corresponding centroids are shown in the following Figures 2-7.



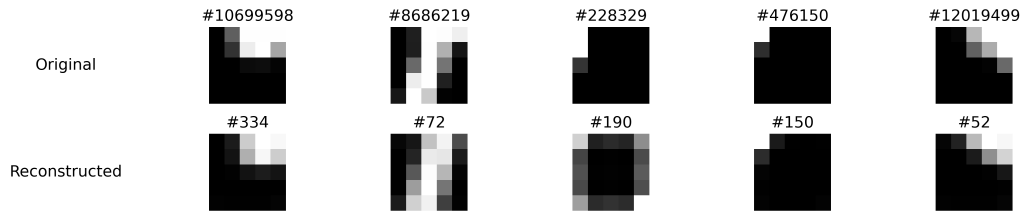Figure 2: The patch and corresponding centroid when $K = 100$



Figure 3: The patch and corresponding centroid when $K = 400$



Figure 4: The patch and corresponding centroid when $K = 1000$



Figure 5: The patch and corresponding centroid when $K = 2000$

Moreover, we plot the Mean Squared Error (MSE) between the original patches and the reconstructed patches (centroids) with different $K$ values, as shown in Figure 8.

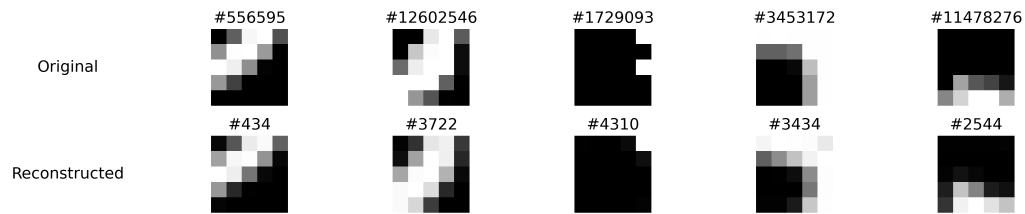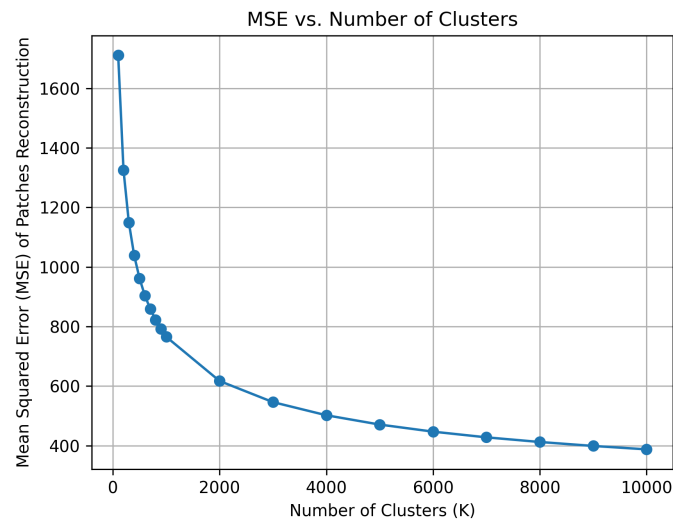With the results above, we can observe that:

Figure 6: The patch and corresponding centroid when $K = 5000$



Figure 7: The patch and corresponding centroid when $K = 10000$



Figure 8: MSE with different clustering number

- Question 1: As K increases, the learned clusters become more and more detailed and similar to the original patches.

- Question 2: With the increase of K, one can reconstruct the patches more accurately. For example, when $K = 100$, the reconstructed patches are not very similar, but when $K = 10000$, the reconstructed patches are very similar to the original patches.

- Question 3: The MSE decreases as K increases, however, the decreasing rate becomes slower when K is large. By elbow method, we can find that the optimal K is around 2000. So we need 2000 clusters to cover the whole patch space, considering the trade-off between the clustering performance and the computational cost.

## 3.2 What are the centroids

In addition to the clustering performance, we also want to understand the meaning of the clusters or centroids. Here we show all the centroids when using $K = 100$ and $K = 1000$ in Figure 9.



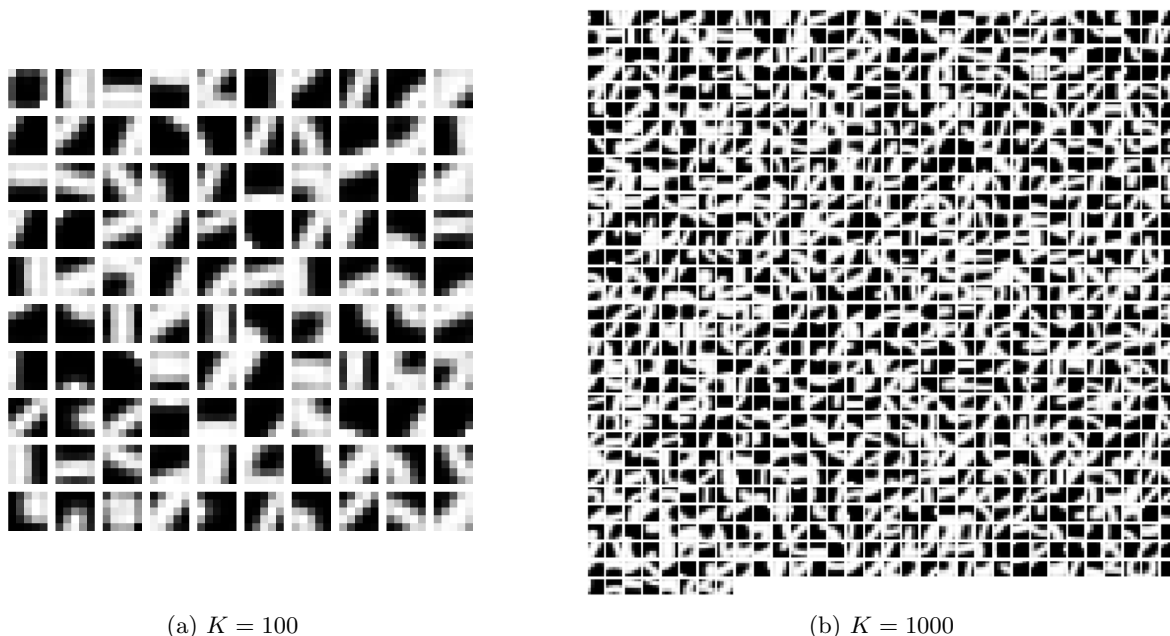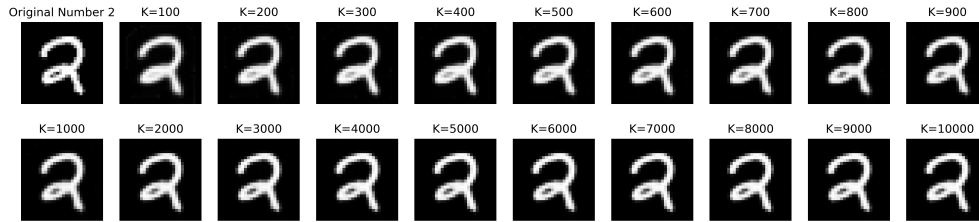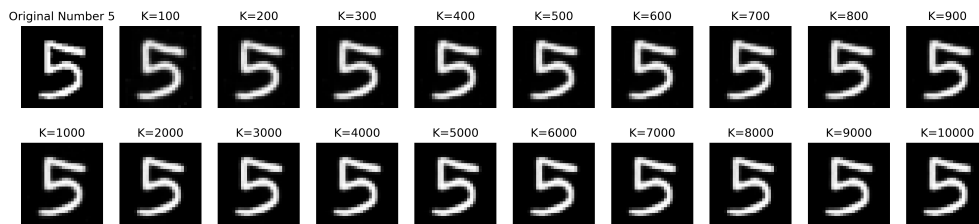(a) $K = 100$        (b) $K = 1000$

Figure 9: Visual representation of cluster centroids for $K = 100$ and $K = 1000$

Question 4: From the results, we can see that the centroids of these clusters are actually the patterns of the handwritten digits. Specifically, they are the edges, corners, and other features of the digits.

## 3.3 Reconstruct a Digit

In this part, we show the reconstruction results of the handwritten digits from both the training set and testing set, by using the learned centroids, as in Figures 10 and 11.

From the results, we can observe that the reconstructed digits are very vague when $K = 100$, but they become more and more clear as $K$ increases. When $K = 10000$, the reconstructed digits are very similar to the original digits. Question 5: This indicates that the digits are constructed from the detailed features of itself, such as the edges, corners, and other features. One digit is just a linear combination of these features / centroids/ clusters, and the more clusters we have, the more accurate the reconstruction will be.

Figure 10: Reconstruction of digit 2 with different $K$



Figure 11: Reconstruction of digit 5 with different $K$

# 4    Conclusion

In this project, we use K-mean clustering to learn the dictionary (cluster) of the $5 \times 5$ patches of the MNIST dataset, and analyze the performance with different K. We can now answer the questions we raised at the beginning:

- When K increases, the learned clusters contain more detailed information on each figure element.

- As shown in Figure 2-7, a $5 \times 5$ MNIST patch can be well reconstructed by the learned clusters. The reconstruction quality is better when K is large.

- Ideally, we want the number of clusters to be as large as possible. However, Figure 8 implies that the marginal return is lower when $K > 2000$, and in practice the training time is huge when K is large. So $K = 2000$ would be a reasonable choice.

- The visualization of these clusters are shown in Figure 9. These clusters are the "elements" of each handwritten digits. They cover some common patterns such as lines, curves, and dots.

- The reconstructed digits are shown in Figure 10-11. With K=100, the reconstructed digits are already recognizable. The reconstruction quality is better with a larger K.

# Appendix

```
1   # https://www.kaggle.com/code/hojjatk/read-mnist-dataset/notebook
2   #
3   # This is a sample Notebook to demonstrate how to read "MNIST Dataset"
4   #
5   import numpy as np # linear algebra
6   import struct
7   from array import array
8
9   #
10  # MNIST Data Loader Class
11  #
12  class MnistDataloader(object):
13      '''
14      MNIST Data Loader
15
16          @type   training_images_filepath: string
17          @param  training_images_filepath: training images file path
18
19          @type   training_labels_filepath: string
20          @param  training_labels_filepath: training labels file path
21
22          @type   test_images_filepath: string
23          @param  test_images_filepath: test images file path
24
25          @type   test_labels_filepath: string
26          @param  test_labels_filepath: test labels file path
27      '''
28      def __init__(self, training_images_filepath, training_labels_filepath,
            test_images_filepath, test_labels_filepath):
29          self.training_images_filepath = training_images_filepath
30          self.training_labels_filepath = training_labels_filepath
31          self.test_images_filepath = test_images_filepath
32          self.test_labels_filepath = test_labels_filepath
33
34      def read_images_labels(self, images_filepath, labels_filepath):
35          '''
36          Read images and labels
37
38                  @type   images_filepath: string
39                  @param  images_filepath: images file path
40
41                  @type   labels_filepath: string
42                  @param  labels_filepath: labels file path
43
44                  @rtype:   (ndarray, ndarray)
45                  @return:  images, labels
46          '''
47          labels = []
48          with open(labels_filepath, 'rb') as file:
49              magic, size = struct.unpack(">II", file.read(8))
50              if magic != 2049:
```

```
51                   raise ValueError('Magic number mismatch, expected 2049, got {}
                        '.format(magic))
52               labels = array("B", file.read())
53
54           with open(images_filepath, 'rb') as file:
55               magic, size, rows, cols = struct.unpack(">IIII", file.read(16))
56               if magic != 2051:
57                   raise ValueError('Magic number mismatch, expected 2051, got {}
                        '.format(magic))
58               image_data = array("B", file.read())
59           images = []
60           for i in range(size):
61               images.append([0] * rows * cols)
62           for i in range(size):
63               img = np.array(image_data[i * rows * cols:(i + 1) * rows * cols])
64               img = img.reshape(28, 28)
65               images[i][:] = img
66
67           return images, labels
68
69       def load_data(self):
70           '''
71           Load MNIST data
72
73               @rtype:    (ndarray, ndarray), (ndarray, ndarray)
74               @return:   (training data, traing labels), (test data, test labels)
75           '''
76           x_train, y_train = self.read_images_labels(self.
                  training_images_filepath, self.training_labels_filepath)
77           x_test, y_test = self.read_images_labels(self.test_images_filepath,
                  self.test_labels_filepath)
78           return (np.array(x_train), np.array(y_train)),(np.array(x_test), np.
                  array(y_test))
```

```
1   from mnist_data_loader import MnistDataloader
2   from os.path   import join
3   import numpy as np
4   from tqdm import tqdm
5
6
7   # Extract 5x5 patches from the 28x28 images
8   def extract_patches(images, patch_size=5, threshold=0):
9       '''
10      Extract patches from images
11
12          @type    images: ndarray
13          @param   images: images
14
15          @type    patch_size: int
16          @param   patch_size: patch size
17
18          @type    threshold: float
19          @param   threshold: default 0 means non-blank patches from the training
```

```
                images

21          @rtype:    ndarray
22          @return:   patches
23      '''
24      patches = []
25      num_pixels = patch_size * patch_size
26      for image in tqdm(images, desc="Extracting patches"):
27          # Slide over the image and extract patches
28          for i in range(image.shape[0] - patch_size + 1):
29              for j in range(image.shape[1] - patch_size + 1):
30                  patch = image[i:i + patch_size, j:j + patch_size]
31                  # Calculate the proportion of non-zero pixels
32                  if np.sum(patch) > 255 * num_pixels * threshold:  # Adjust the
                        threshold as needed
33                      patches.append(patch.flatten())
34      return np.array(patches)



37  # Extract 5x5 patches from the 28x28 images
38  def extract_nonblank_patches_from_one_image(image, patch_size=5, threshold=0):
39      '''
40      Extract nonblank patches from one image

42          @type    image: ndarray
43          @param   image: image

45          @type    patch_size: int
46          @param   patch_size: patch size

48          @type    threshold: float
49          @param   threshold: default 0 means non-blank patches from the training
                images

51          @rtype:    ndarray
52          @return:   patches
53      '''
54      patches = []
55      num_pixels = patch_size * patch_size
56      # Slide over the image and extract patches
57      for i in range(image.shape[0] - patch_size + 1):
58          for j in range(image.shape[1] - patch_size + 1):
59              patch = image[i:i + patch_size, j:j + patch_size]
60              # Calculate the proportion of non-zero pixels
61              if np.sum(patch) > 255 * num_pixels * threshold:  # Adjust the
                    threshold as needed
62                  patches.append(patch.flatten())
63      return np.array(patches)



66  def extract_all_patches_from_one_image(image, patch_size=5):
67      '''
68      Extract all patches (blank and nonblank) from one image

```

```
70          @type    image: ndarray
71          @param   image: image
72
73          @type    patch_size: int
74          @param   patch_size: patch size
75
76          @rtype:    ndarray
77          @return:   patches
78      '''
79      patches = []
80      # Slide over the image and extract patches
81      for i in range(image.shape[0] - patch_size + 1):
82          for j in range(image.shape[1] - patch_size + 1):
83              patch = image[i:i + patch_size, j:j + patch_size]
84              patches.append(patch.flatten())
85      return np.array(patches)
86
87
88
89  if __name__ == "__main__":
90      # Set file paths based on added MNIST Datasets
91      input_path = 'MNIST_ORG/'
92      training_images_filepath = join(input_path, 'train-images.idx3-ubyte')
93      training_labels_filepath = join(input_path, 'train-labels.idx1-ubyte')
94      test_images_filepath = join(input_path, 't10k-images.idx3-ubyte')
95      test_labels_filepath = join(input_path, 't10k-labels.idx1-ubyte')
96
97      # Load MINST dataset
98      mnist_dataloader = MnistDataloader(training_images_filepath,
              training_labels_filepath, test_images_filepath, test_labels_filepath)
99      (x_train, y_train), (x_test, y_test) = mnist_dataloader.load_data()
100
101
102     # Extract non-blank patches from the training data
103     patches = extract_patches(x_train)
104     np.save("patches.npy", patches)
```

```
1  import torch
2  from tqdm import tqdm
3
4  def kmeans_pytorch(X, n_clusters, n_iters=100, tolerance=1e-4):
5      """
6      Performs k-means clustering using PyTorch.
7
8      Parameters:
9          X (torch.Tensor): The input data, a tensor of shape (n_samples,
              n_features).
10         n_clusters (int): The number of clusters to form.
11         n_iters (int): Maximum number of iterations of the k-means algorithm.
12         tolerance (float): Tolerance to declare convergence.
13
14     Returns:
15         centers (torch.Tensor): Cluster centers, a tensor of shape (n_clusters
```

```
                   , n_features).
16             labels (torch.Tensor): Index of the cluster each sample belongs to.
17         """
18         # Randomly choose cluster centers from the input data at the start.
19         indices = torch.randperm(X.size(0))[:n_clusters]
20         centers = X[indices]
21
22         for _ in tqdm(range(n_iters), desc="K-means"):
23             # Compute distances from data points to the centroids
24             distances = torch.cdist(X, centers)
25             # Assign clusters
26             labels = torch.argmin(distances, dim=1)
27             # Compute new centers
28             new_centers = torch.stack([X[labels == i].mean(dim=0) for i in range(
                   n_clusters)])
29
30             # Check for convergence
31             if torch.norm(centers - new_centers) < tolerance:
32                 break
33
34             centers = new_centers
35
36         return centers, labels
37
38
39 if __name__ == "__main__":
40     # Example usage
41     # Creating some data
42     torch.manual_seed(0)
43     data = torch.randn(100, 2)  # 100 data points, 2 dimensions
44
45     # Clustering
46     centers, labels = kmeans_pytorch(data, n_clusters=3)
47     print("Cluster_centers:\n", centers)
48     print("Cluster_labels:\n", labels)
```

```
1  import numpy as np
2  import matplotlib.pyplot as plt
3  from sklearn.cluster import KMeans
4  import joblib
5  from tqdm import tqdm
6
7  # Load patches
8  try:
9      patches = np.load("K-means/Code/patches.npy")
10 except FileNotFoundError:
11     print("cd_to_folder_EEC289A,_run_extract_patches.py_first.")
12     exit()
13
14 K = [100, 200, 300, 400, 500, 600, 700, 800, 900, 1000, 2000, 3000, 4000,
       5000, 6000, 7000, 8000, 9000, 10000]
15
16 # Perform K-means clustering
```

```
17  for n_clusters in tqdm(K, desc="Clustering"):
18      kmeans = KMeans(n_clusters=n_clusters, random_state=0).fit(patches)
19
20      # Save the model
21      joblib.dump(kmeans, f"K-means/Result/Model/{n_clusters}-clusters-model.
            joblib")
22
23
24  # # Load the model
25  # model = joblib.load("../Result/Model/100-clusters-model.joblib")
```

```
1   import numpy as np
2   import matplotlib.pyplot as plt
3   from sklearn.cluster import KMeans
4   from sklearn.metrics import mean_squared_error
5   import joblib
6
7
8   def visualize_reconstruction(n_clusters, patches, model):
9       '''
10      Visualize the original and reconstructed patches for a few random samples.
11
12          @type    n_clusters: int
13          @param   n_clusters: K
14
15          @type    patches: ndarray
16          @param   patches: patches
17
18          @type    model: sklearn model
19          @param   model: the fitted sklearn KMeans model
20      '''
21      # Predict the cluster for each patch
22      labels = model.labels_
23
24      # Get the cluster centers
25      centroids = model.cluster_centers_
26
27      # Pick random patches for display
28      num_samples = 5  # Number of random samples to pick
29      indices = np.random.choice(range(len(patches)), num_samples, replace=False
            )
30
31      # Plotting the original and reconstructed patches
32      fig, axs = plt.subplots(2, num_samples+1, figsize=(15, 3))  # 2 rows:
            originals and reconstructions
33
34      # Set labels for the rows
35      axs[0, 0].text(0.5, 0.5, 'Original', verticalalignment='center',
            horizontalalignment='center', transform=axs[0, 0].transAxes, fontsize
            = 15)
36      axs[1, 0].text(0.5, 0.5, 'Reconstructed', verticalalignment='center',
            horizontalalignment='center', transform=axs[1, 0].transAxes, fontsize
            = 15)
```

```python
37        axs[0, 0].axis('off')
38        axs[1, 0].axis('off')
39
40        for i, idx in enumerate(indices):
41            i += 1  # Adjust index for the extra label column
42
43            # Original patches
44            axs[0, i].imshow(patches[idx].reshape(5, 5), cmap='gray')
45            axs[0, i].axis('off')
46            axs[0, i].set_title('#{}'.format(idx), fontsize = 15)
47
48            # Reconstructed patches
49            reconstructed_patch = centroids[labels[idx]].reshape(5, 5)
50            axs[1, i].imshow(reconstructed_patch, cmap='gray')
51            axs[1, i].axis('off')
52            axs[1, i].set_title('#{}'.format(labels[idx]), fontsize = 15)
53
54        plt.tight_layout()
55        # plt.show()
56        plt.savefig(f"K-means/Result/Reconstruction/{n_clusters}-clusters-
              reconstruction.png", dpi=300)
57        plt.close()
58
59
60    def mse_reconstruction(patches, model):
61        '''
62        Calculate the mean squared error between the original and reconstructed
              patches.
63
64            @type    patches: ndarray
65            @param   patches: patches
66
67            @type    model: sklearn model
68            @param   model: the fitted sklearn KMeans model
69        '''
70        # Predict the cluster for each patch
71        labels = model.labels_
72
73        # Get the cluster centers
74        centroids = model.cluster_centers_
75
76        # Calculate the mean squared error
77        reconstruction = centroids[labels]
78        mse = mean_squared_error(patches, reconstruction)
79
80        return mse
81
82
83
84
85
86    if __name__ == "__main__":
87        K = [100, 200, 300, 400, 500, 600, 700, 800, 900, 1000, 2000, 3000, 4000,
              5000, 6000, 7000, 8000, 9000, 10000]
```

```
88
89        # Load patches
90        try:
91             patches = np.load("K-means/Code/patches.npy")
92        except FileNotFoundError:
93             print("cd_to_folder_EEC289A,_run_extract_patches.py_first.")
94             exit()
95
96        mse_K = []
97
98        for n_clusters in K:
99             # Load the model
100            try:
101                 model = joblib.load(f"K-means/Result/Model/{n_clusters}-clusters-
                        model.joblib")
102            except FileNotFoundError:
103                 print(f"cd_to_folder_EEC289A,_run_run_kmeans.py_first.")
104                 exit()
105
106            visualize_reconstruction(n_clusters, patches, model)
107
108            mse = mse_reconstruction(patches, model)
109            mse_K.append(mse)
110
111        # Plot the mean squared error
112        plt.plot(K, mse_K, marker='o')
113        plt.xlabel('Number_of_Clusters_(K)')
114        plt.ylabel('Mean_Squared_Error_(MSE)_of_Patches_Reconstruction')
115        plt.title('MSE_vs._Number_of_Clusters')
116        plt.grid(True)
117        plt.savefig("K-means/Result/Reconstruction/MSE_vs_K.png", dpi=300)
```

```
1    import numpy as np
2    import matplotlib.pyplot as plt
3    from sklearn.cluster import KMeans
4    import joblib
5    from os.path import join
6
7    from mnist_data_loader import MnistDataloader
8    from extract_patches import extract_all_patches_from_one_image,
         extract_nonblank_patches_from_one_image
9
10
11
12   def calculate_positions(image_shape, patch_size):
13        '''
14        Calculate the positions of all patches in an image
15
16            @type    image_shape: tuple
17            @param   image_shape: image shape
18
19            @type    patch_size: int
20            @param   patch_size: patch size
```

```
21
22          @rtype:      list
23          @return:   positions
24      '''
25      positions = []
26      for i in range(image_shape[0] - patch_size + 1):
27          for j in range(image_shape[1] - patch_size + 1):
28              positions.append((i, j))
29      return positions
30
31
32
33  def reconstruct_digit(digit_image, model, patch_size = 5):
34      '''
35      Reconstruct a digit image using a KMeans model
36
37          @type    digit_image: ndarray
38          @param   digit_image: digit image
39
40          @type    model: sklearn model
41          @param   model: KMeans model
42
43          @type    patch_size: int
44          @param   patch_size: patch size
45
46          @rtype:    ndarray
47          @return:   reconstructed image
48      '''
49      # Extract all patches and calculate positions
50      all_patches = extract_all_patches_from_one_image(digit_image)
51      positions = calculate_positions(digit_image.shape, patch_size)
52
53      # Extract non-blank patches
54      nonblank_patches = extract_nonblank_patches_from_one_image(digit_image)
55
56      # Map nonblank patches to their positions
57      nonblank_indices = [i for i, patch in enumerate(all_patches) if np.sum(
              patch) > 0]
58
59      # Predict clusters for non-blank patches
60      labels = model.predict(nonblank_patches)
61
62      # Initialize the reconstructed image with zeros
63      reconstructed_image = np.zeros_like(digit_image, dtype=float)
64      count_matrix = np.zeros_like(digit_image, dtype=float)
65
66      # Add centroids to the corresponding positions
67      centroids = model.cluster_centers_
68      for label, idx in zip(labels, nonblank_indices):
69          i, j = positions[idx]
70          reconstructed_image[i:i+patch_size, j:j+patch_size] += centroids[label
              ].reshape(patch_size, patch_size)
71          count_matrix[i:i+patch_size, j:j+patch_size] += 1
72
```

```
73        # Avoid division by zero
74        count_matrix[count_matrix == 0] = 1
75        reconstructed_image /= count_matrix
76
77        return reconstructed_image
78
79
80
81
82
83    if __name__ == "__main__":
84        # Set file paths based on added MNIST Datasets
85        input_path = 'K-means/Code/MNIST_ORG/'
86        training_images_filepath = join(input_path, 'train-images.idx3-ubyte')
87        training_labels_filepath = join(input_path, 'train-labels.idx1-ubyte')
88        test_images_filepath = join(input_path, 't10k-images.idx3-ubyte')
89        test_labels_filepath = join(input_path, 't10k-labels.idx1-ubyte')
90
91
92        # Load MINST dataset
93        mnist_dataloader = MnistDataloader(training_images_filepath,
              training_labels_filepath, test_images_filepath, test_labels_filepath)
94        (x_train, y_train), (x_test, y_test) = mnist_dataloader.load_data()
95
96
97        # Example of reconstructing a digit
98        digit_idx = np.random.randint(0, len(x_test))
99        digit_image = x_test[digit_idx]
100
101
102
103       K = [100, 200, 300, 400, 500, 600, 700, 800, 900, 1000, 2000, 3000, 4000,
              5000, 6000, 7000, 8000, 9000, 10000]
104
105
106       # Prepare the figure for subplots
107       fig, axs = plt.subplots(2, 10, figsize=(18, 4))  # 1 row, columns for each
              K plus one for the original
108
109       # Display the original digit in the first column
110       axs[0, 0].imshow(digit_image, cmap='gray')
111       axs[0, 0].set_title('Original Number {}'.format(y_train[digit_idx]))
112       axs[0, 0].axis('off')
113
114
115
116       for i_fig, n_clusters in enumerate(K):
117           # Load the model
118           try:
119               model = joblib.load(f"K-means/Result/Model/{n_clusters}-clusters-
                      model.joblib")
120           except FileNotFoundError:
121               print(f"cd to folder EEC289A, run run_kmeans.py first.")
122               exit()
```

17

```
123
124            # Reconstruct the digit
125            reconstructed_image = reconstruct_digit(digit_image, model)
126
127            # Display reconstructed digit for this K
128            if i_fig < 9:
129                axs[0, i_fig + 1].imshow(reconstructed_image, cmap='gray')
130                axs[0, i_fig + 1].set_title(f'K={n_clusters}')
131                axs[0, i_fig + 1].axis('off')
132            else:
133                axs[1, i_fig - 9].imshow(reconstructed_image, cmap='gray')
134                axs[1, i_fig - 9].set_title(f'K={n_clusters}')
135                axs[1, i_fig - 9].axis('off')
136
137        # plt.tight_layout()
138        plt.savefig(f"K-means/Result/Digits/reconstruct-test-digit.png", dpi=300)
```