**Name (netid):** Your Name (Your Netid)
**CS 441 - HW 4: Dealing with Data**

Complete the claimed points and sections below.

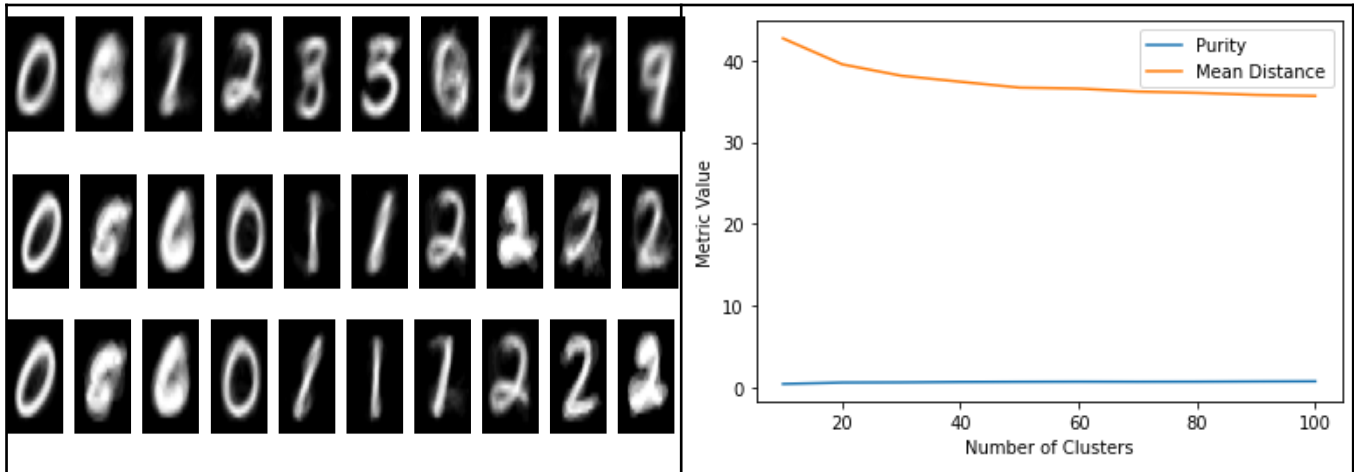**Total Points Claimed**                                       **[ ] / 142**

1. Clustering and Fast Retrieval
   a. Test Kmeans Purity & Centroids     [ ] / 15
   b. Questions     [ ] / 10
   c. Fast 1-NN Retrieval     [ ] / 15
2. Estimating PDFs
   a. Histograms     [ ] / 10
   b. Clustering     [ ] / 10
   c. Gaussian Mixture Model     [ ] / 15
3. PCA and Data Compression
   a. Display Principal Components     [ ] / 5
   b. Scatter Plot     [ ] / 5
   c. Plot cumulative explained variance     [ ] / 5
   d. Time & Accuracy     [ ] / 10
4. Stretch Goals
   a. Rotate Using PCA and comparison
      To original approach     [ ] / 15
   b. Try Part 2 with your own images     [ ] / 10
   c. Plot Using t-SNE and MDS     [ ] / 15
   d. Completed HW3 survey by DATE     [ ] / 2

## 1. Clustering & Fast Retrieval
   a. Test KMeans Purity
      [15]

| K Vs Purity Plot | K Vs Mean_Distance Plot |
| --- | --- |
| | |

Paste images of centroids for K = 10, K = 20, and K=30 below (three rows).

b. Questions                                                                 [10]

    i.   As you increase K, do you expect the purity to increase?  Why or why not?

**Increasing K in k-means clustering may or may not increase the purity of the clusters. It depends on the data and the structure of the clusters. If the clusters are well-separated and distinct, increasing K may improve the purity as it allows the algorithm to capture more subtle variations in the data. However, if the clusters are overlapping or poorly separated, increasing K may lead to fragmentation of the clusters and decrease the purity.**

    ii.   In a given run, is the average distance of a sample to centroid guaranteed to monotonically decrease  with each iteration (i.e. cannot increase)? Why or why not?

**No, the average distance of a sample to the centroid is not guaranteed to monotonically decrease with each iteration, as k-means is sensitive to the initial centroid positions and can converge to local optima.**

    iii.   If you do enough iterations, is Kmeans guaranteed to give you the optimal clustering that minimizes the sum of distances between each sample and its center? Why or why not?

**No, Kmeans is not guaranteed to give the optimal clustering that minimizes the sum of distances between each sample and its center, even after enough iterations, because it can get stuck in local optima. The algorithm's final result depends on the initial choice of centroids and can converge to a suboptimal solution.**

    iv.   Does improving the Kmeans objective (i.e. achieving lower mean squared error) necessarily improve expected purity ? Why or why not?

**Improving the Kmeans objective does not necessarily improve the expected purity, because the objective only focuses on minimizing the sum of distances between each sample and its center, without taking into account the distribution of the samples across the clusters. Therefore, it is possible to have a low objective value, but with clusters that have a high degree of overlap and low purity.**

c. Fast Retrievals

    i.   Brute Force:                                                     [5]

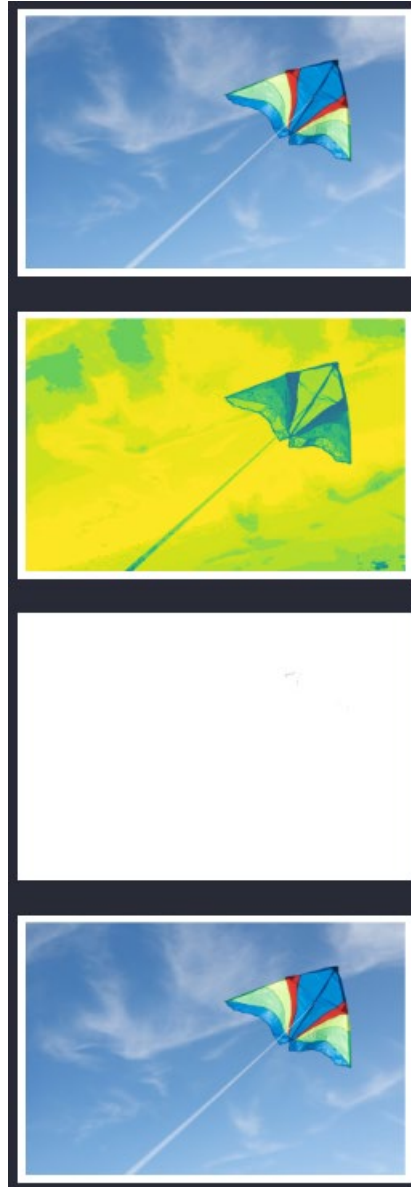| Test Error | Time to Add | Time to Search |
|---|---|---|
| 0.0089 | 0.1147sZ | 4.2405s |

    ii.   LSH:                                                             [5]

| Test Error | Time to Add | Time to Search | Nbits parameter |
|---|---|---|---|
| 0.3943 | 0.1554s | 0.5845s | 8 |

## 2. Estimating PDFs

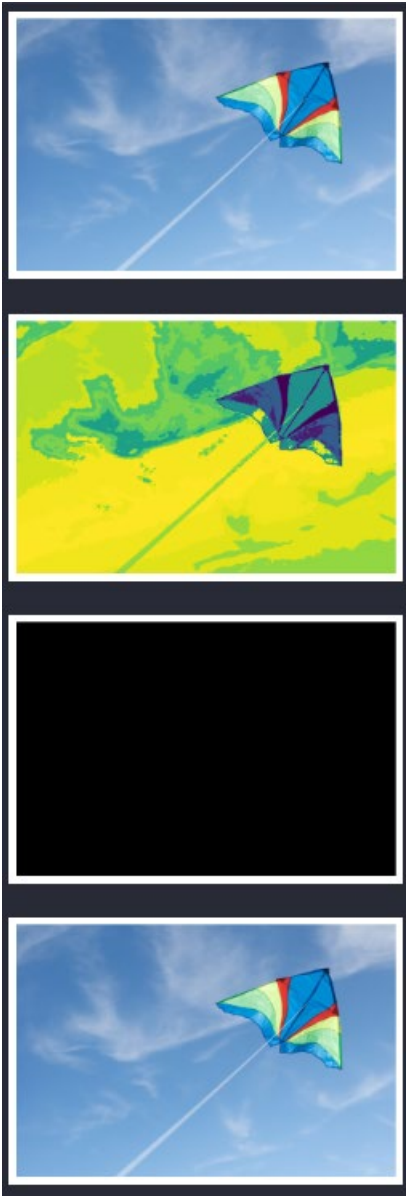Include the generated images (score map, thresholded score map, thresholded RGB) from the display code.

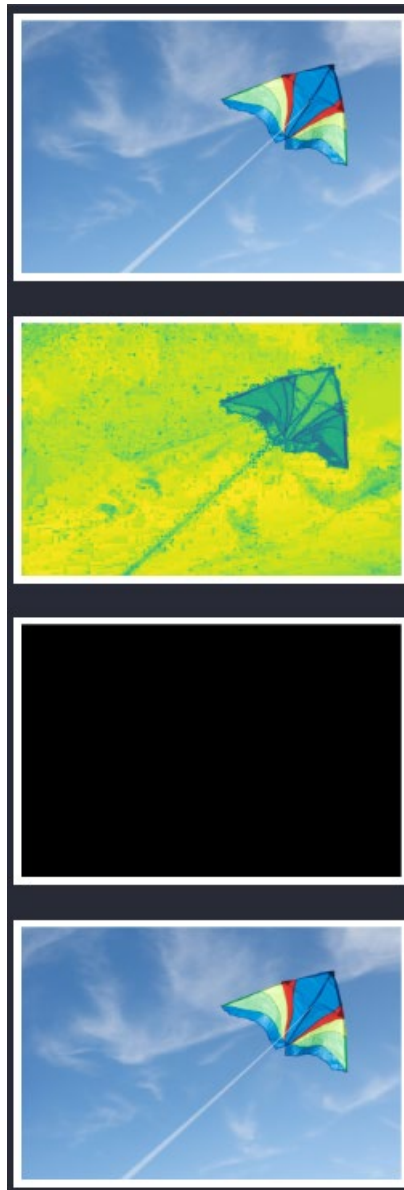a. Histogram: [10]
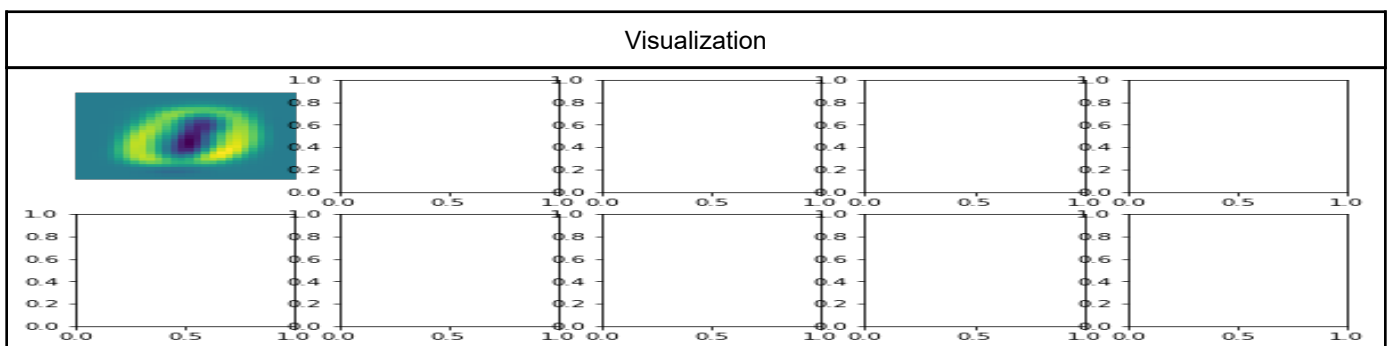


b. Clustering: [10]

c. Gaussian mixture Model: [15]

## 3. PCA and Data Compression

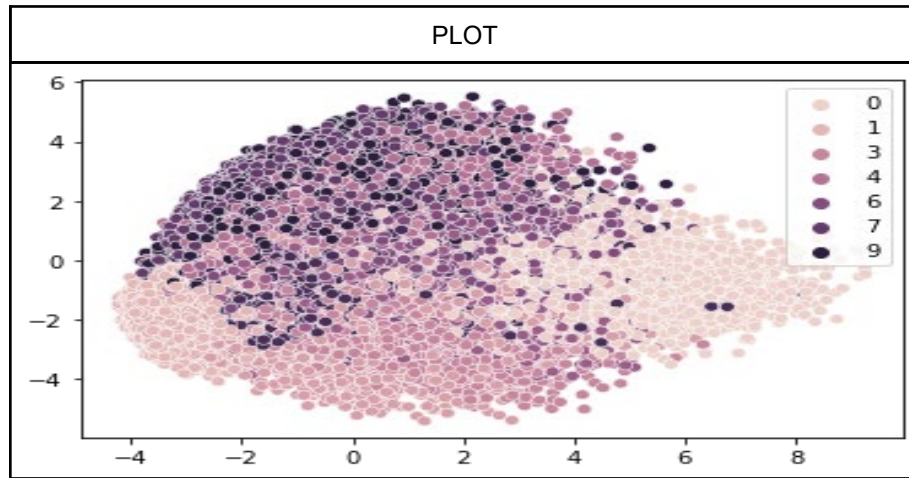1. First 10 principal components                                    [5]
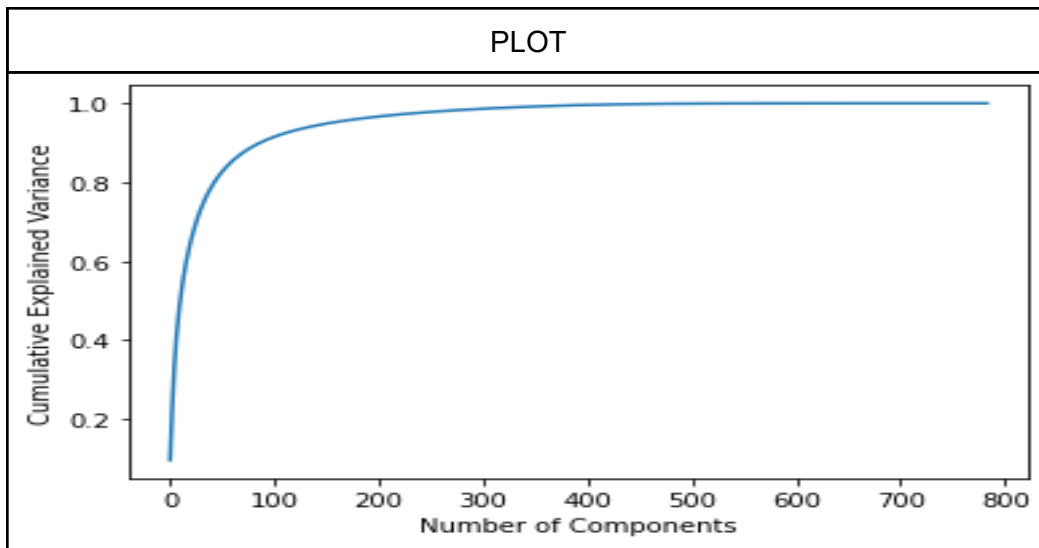
2. Scatterplot                                                                    [5]



3. Cumulative explained Variance                                                   [5]

PLOT

4. Faiss                                                                                    [10]

|  | Total Time | Test Error | Dimensions |
|---|---|---|---|
| Brute Force (PCA) | 20.55 | 0.0331 | 784 |
| Brute Force | 0.1147 | 0.0089 | 784 |
| LSH | 0.1554 | 0.3943 | 784 |

Note: the last two rows are copied from 1.c for reference.

## 4. Stretch Goals
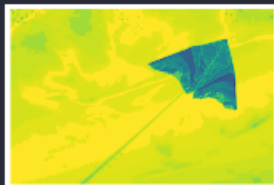a. PDFS after using PCA to rotate your data                                                 [15]

## Q4-P1

```python
1  # set number of bins K
2  K = 16
3  # reshape pixels to (h*w, 3)
4  im_pixels = im.reshape(-1, 3)
5
6  # center the data
7  im_pixels_centered = im_pixels - np.mean(im_pixels, axis=0)
8
9  # perform PCA
10 covariance = np.cov(im_pixels_centered.T)
11 eigenvalues, eigenvectors = np.linalg.eig(covariance)
12 # sort eigenvectors in descending order of eigenvalues
13 idx = eigenvalues.argsort()[::-1]
14 eigenvectors = eigenvectors[:, idx]
15
16 # project data onto principal components
17 im_pixels_rotated = np.dot(im_pixels_centered, eigenvectors)
18
19 # convert continuous values to discrete values ranging from 0 to K-1
20 # e.g. x-->min(int(x*K), K-1) if x ranges from 0 to 1
21 discrete_im_pixels = np.minimum((im_pixels_rotated*K).astype(int), K-1)
22
23 # get pdf for each dimension using estimate_discrete_pdf
24 pdfs = [estimate_discrete_pdf(discrete_im_pixels[:, i], K) for i in range(3)]
25
26 # estimate score for each pixel in full image according to log pdfs
27 scores = np.zeros(im.shape[:2])
28 for i in range(im.shape[0]):
29     for j in range(im.shape[1]):
30         scores[i, j] = np.sum([np.log(pdfs[k][discrete_im_pixels[i*im.shape[1]+j, k]]) for k in range(3)])
31
32 # display
33 display_score_maps(im, scores, -20)
34
✓ 4.2s
```
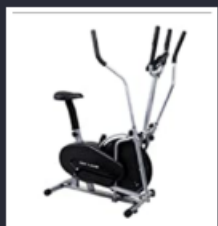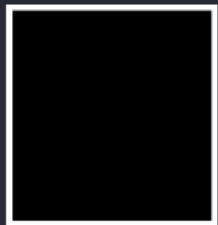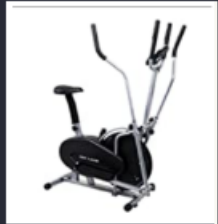


b. Apply Part 2 to your own choice of image, with the same deliverables          [10]

## Method 1 (Per-dimension discrete pdf)

```python
1  # set number of bins K
2  K = 16
3  # reshape pixels to (h*w, 3)
4  im_pixels = im.reshape(-1, 3)
5  # convert continuous values to discrete values ranging from 0 to K-1
6  # e.g. x-->min(int(x*K), K-1) if x ranges from 0 to 1
7  discrete_im_pixels = np.minimum((im_pixels*K).astype(int), K-1)
8  # get pdf for each dimension using estimate_discrete_pdf
9  pdfs = [estimate_discrete_pdf(discrete_im_pixels[:, i], K) for i in range(3)]
10 # estimate score for each pixel in full image according to log pdfs
11 scores = np.zeros(im.shape[:2])
12 for i in range(im.shape[0]):
13     for j in range(im.shape[1]):
14         scores[i, j] = np.sum([np.log(pdfs[k][discrete_im_pixels[i*im.shape[1]+j, k]]) for k in range(3)])
15
16 # display
17 display_score_maps(im, scores, -20)
```
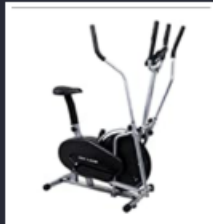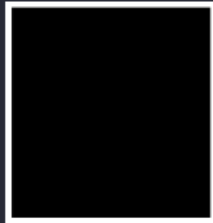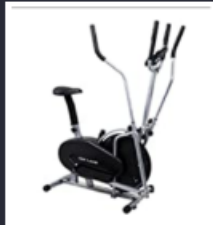
✓ 0.7s

## Method 2 (K-means)

```python
import faiss

# set K
K = 16
# reshape pixels to (h*w, 3)
im_pixels = im.reshape(-1, 3)
# discretize all three color channels together using KMeans
kmeans = faiss.Kmeans(3, K)
kmeans.train(im_pixels.astype(np.float32))
discrete_im_pixels = kmeans.index.search(im_pixels.astype(np.float32), 1)[1][:,0]
# get pdf over discrete values
pdf = estimate_discrete_pdf(discrete_im_pixels, K)
# estimate score for each pixel in full image according to log pdfs
scores = np.zeros(im.shape[:2])
for i in range(im.shape[0]):
    for j in range(im.shape[1]):
        scores[i, j] = np.log(pdf[discrete_im_pixels[i*im.shape[1]+j]])
# display
display_score_maps(im, scores, -20)
```
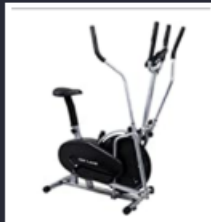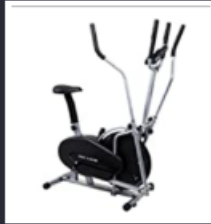
✓ 0.3s

## Method 3 (GMM)

```python
1  from sklearn.mixture import GaussianMixture
2
3  # reshape pixels to (h*w, 3)
4  im_pixels = im.reshape(-1, 3)
5  # get joint pdf using GMMs (choose number of components and other parameters)
6  gmm = GaussianMixture(n_components=10)
7  gmm.fit(im_pixels)
8  pdf = gmm.score_samples(im_pixels)
9  # estimate score for each pixel in full image according to log pdfs
10 scores = np.reshape(pdf, im.shape[:2])
11 # display
12 display_score_maps(im, scores, -20)
```
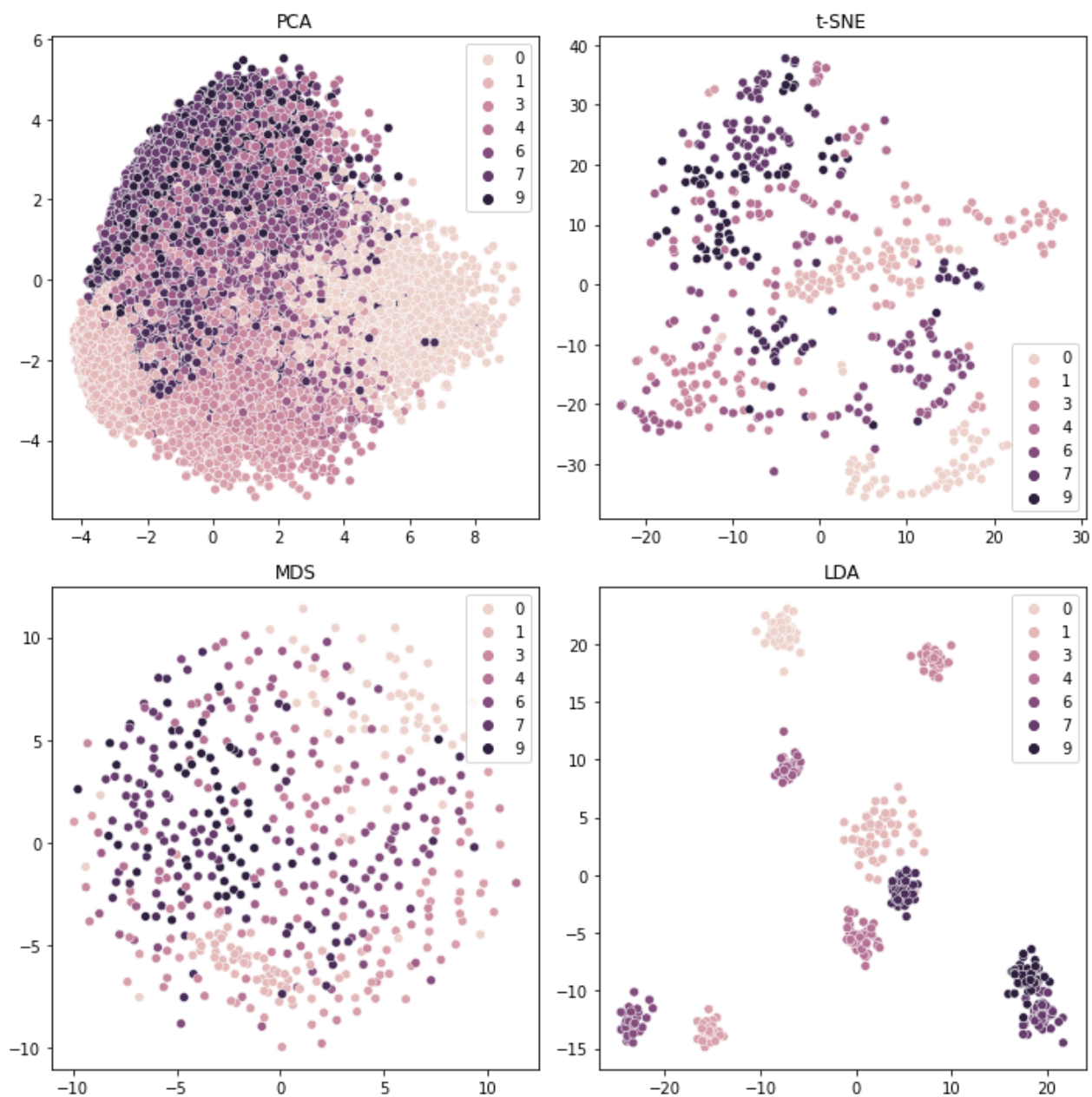✓ 3.0s









c.  Scatterplots using at least two of t-SNE, MDS, and Linear Discriminant Analysis  [15]

```
1    # Import the required libraries
2    from sklearn.manifold import TSNE, MDS
3    from sklearn.discriminant_analysis import LinearDiscriminantAnalysis as LDA
4
5    # Transform the data using t-SNE
6    tsne = TSNE(n_components=2, random_state=42)
7    x_train_tsne = tsne.fit_transform(x_train[train_indices['s']])
8
9    # Transform the data using MDS
10   mds = MDS(n_components=2, random_state=42)
11   x_train_mds = mds.fit_transform(x_train[train_indices['s']])
12
13   # Transform the data using LDA
14   lda = LDA(n_components=2)
15   x_train_lda = lda.fit_transform(x_train[train_indices['s']], y_train[train_indices['s']])
16
17   # Display the scatter plots
18   fig, axs = plt.subplots(2, 2, figsize=(10, 10))
19
20   # Plot PCA scatterplot
21   sns.scatterplot(x_train_pca_2d[:,0], x_train_pca_2d[:,1], hue=y_train, ax=axs[0,0])
22   axs[0,0].set_title('PCA')
23
24   # Plot t-SNE scatterplot
25   sns.scatterplot(x_train_tsne[:,0], x_train_tsne[:,1], hue=y_train[train_indices['s']], ax=axs[0,1])
26   axs[0,1].set_title('t-SNE')
27
28   # Plot MDS scatterplot
29   sns.scatterplot(x_train_mds[:,0], x_train_mds[:,1], hue=y_train[train_indices['s']], ax=axs[1,0])
30   axs[1,0].set_title('MDS')
31
32   # Plot LDA scatterplot
33   sns.scatterplot(x_train_lda[:,0], x_train_lda[:,1], hue=y_train[train_indices['s']], ax=axs[1,1])
34   axs[1,1].set_title('LDA')
35
36   plt.tight_layout()
37   plt.show()
```

**Acknowledgments / Attribution**

"None".