

## Problem 1

In this problem, we will look at image compression using SVD, following the lines of the well-known "Eigenfaces" experiment. The basic concept is to represent an image (in grayscale) of size  $m \times n$  as an  $m \times n$  real matrix  $M$ . SVD is then applied to this matrix to obtain  $U$ ,  $S$ , and  $V^T$  such that  $M = USV^T$ . Here  $U$  and  $V$  are the matrices whose columns are the left and right singular vectors, respectively, and  $S$  is a diagonal  $m \times n$  matrix consisting of the singular values of  $M$ . The number of non-zero singular values is the rank of  $M$ . By using just the largest  $k$  singular values (and corresponding left and right singular vectors), one obtains the best rank- $k$  approximation to  $M$ .

The following code returns the dataset of 400 images.

```
1 import matplotlib.pyplot as plt
2 import numpy as np
3 from numpy.linalg import svd
4 from sklearn import datasets
5
6 data = datasets.fetch_olivetti_faces()
7 images = data.images
```

(a)

Given an  $m \times n$  image  $M$  and its rank- $k$  approximation  $A$ , we can measure the reconstruction error using mean  $\ell_1$  error:

$$\text{error}_{\ell_1}(M, A) = \frac{1}{mn} \|M - A\|_1 = \frac{1}{mn} \sum_{i=1}^m \sum_{j=1}^n |M_{i,j} - A_{i,j}|.$$

For  $k = 1, \dots, 30$ , take the average rank- $k$  reconstruction error over all images in the dataset, and plot a curve of average reconstruction error as a function of  $k$ .

```
1 # Returns the best rank-k approximation to M
2 def svd_reconstruct(M, k):
3     U, S, Vt = np.linalg.svd(M, full_matrices=False)
4     S_k = np.diag(S[:k])
5     M_k = U[:, :k] @ S_k @ Vt[:k, :]
6     return M_k
7
8 def mean_l1_error(M, A):
9     return np.mean(np.abs(M - A))
10
11 max_k = 30
12 errors = np.zeros(max_k)
13
14 for k in range(1, max_k + 1):
15     total_error = 0
16     for image in images:
17         M = image
18         A_k = svd_reconstruct(M, k)
19         total_error += mean_l1_error(M, A_k)
20     errors[k-1] = total_error / len(images)
21
22 plt.plot(range(1, max_k + 1), errors, marker='o')
23 plt.xlabel('k')
24 plt.ylabel('Average Reconstruction Error (Mean l1)')
25 plt.title('Average Rank-k Reconstruction Error')
26 plt.grid(True)
27 plt.show()
```

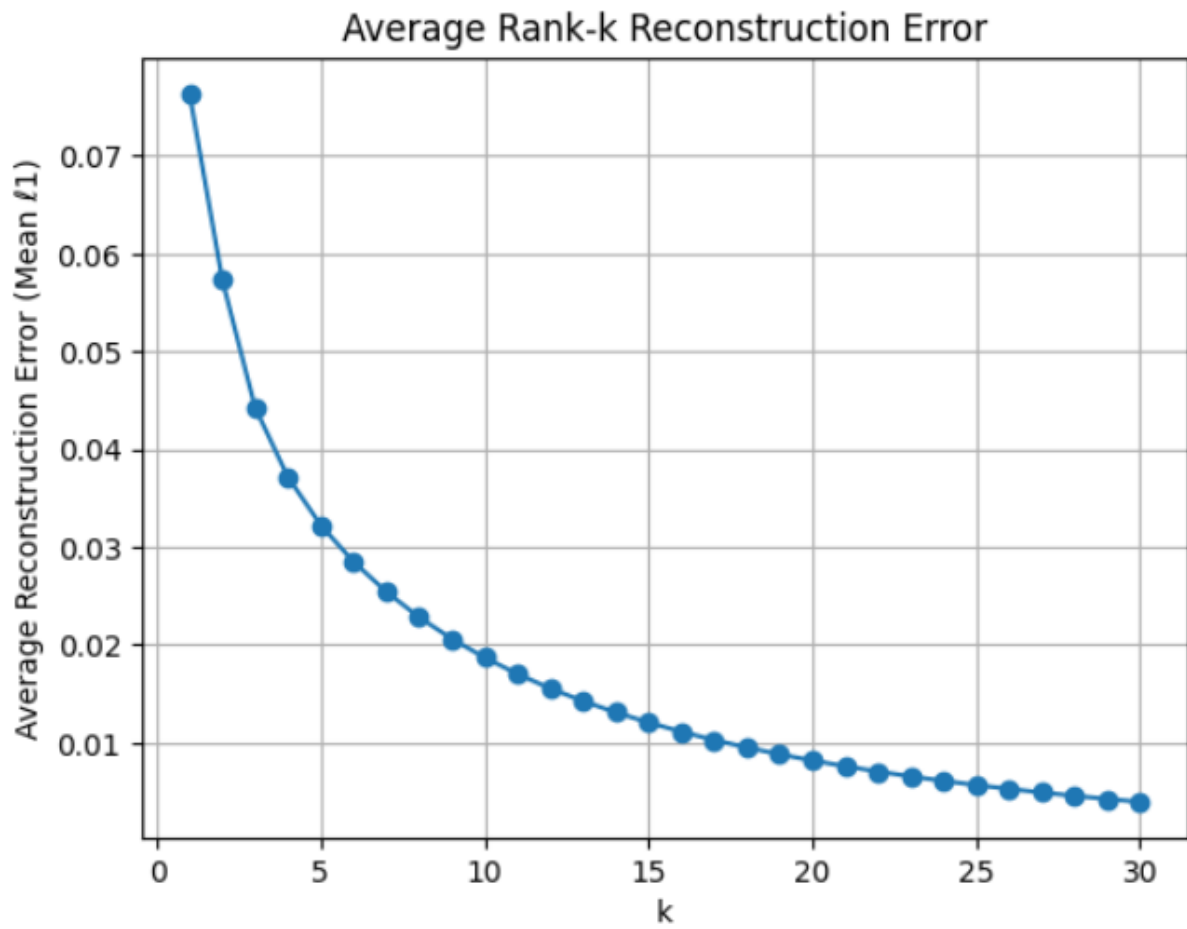


Figure 1: k-rank Reconstruction Error

(b)

Pick any image in the dataset, and display the following side-by-side as images: the original, and the best rank- $k$  approximations for  $k = 10, 20, 30, 40$ . You will find the `imshow` method in `matplotlib` useful for this; pass in `cmap='gray'` to render in grayscale. Feel free to play around further.

```

1 # Pick any image from the dataset (e.g., the first image)
2 image = images[5]
3
4 # Compute rank-k approximations
5 k_values = [10, 20, 30, 40]
6 approximations = [svd_reconstruct(image, k) for k in k_values]
7
8 # Plot the original and approximated images
9 plt.figure(figsize=(10, 5))
10
11 # Original Image
12 plt.subplot(1, len(k_values) + 1, 1)
13 plt.imshow(image, cmap='gray')
14 plt.title('Original')
15 plt.axis('off')
16
17 # Approximations
18 for i, (k, approx) in enumerate(zip(k_values, approximations)):
19     plt.subplot(1, len(k_values) + 1, i + 2)
20     plt.imshow(approx, cmap='gray')
21     plt.title(f'Rank-{k}')
22     plt.axis('off')
23

```

```
plt.show()
```



Figure 2: k-rank Approximations

## Problem 2

In this problem, we visualize the Wisconsin breast cancer dataset in two dimensions using PCA.

First, rescale the data so that every feature has a mean of 0 and a standard deviation of 1 across the various points in the dataset. You may find `sklearn.preprocessing.StandardScaler` useful for this.

Next, compute the top two principal components of the dataset using PCA, and for every data point, compute its coordinates (i.e., projections) along these two principal components. You should do this in two ways:

1. By using SVD directly. Do not use any PCA built-ins.
2. By using `sklearn.decomposition.PCA`.

The two approaches should give exactly the same result, and this also acts as a check. (But note that the signs of the singular vectors may be flipped in the two approaches since singular vectors are only determined uniquely up to sign. If this happens, flip signs to make everything identical again.)

Your final goal is to make a scatterplot of the dataset in 2 dimensions, where the x-axis is the first principal component and the y-axis is the second. Color the points by their diagnosis (malignant or benign). Do this for both approaches. Your plots should be identical. Does the data look roughly separable already in 2 dimensions?

```
1 from sklearn.preprocessing import StandardScaler
2 from sklearn.decomposition import PCA
3 cancer = datasets.load_breast_cancer()
4
5 X = cancer.data # feature matrix
6 y = cancer.target # labels (0: malignant, 1: benign)
7
8 # Standardize the data
9 scaler = StandardScaler()
10 X_scaled = scaler.fit_transform(X)
11
12 # Perform SVD
13 U, S, Vt = np.linalg.svd(X_scaled, full_matrices=False)
14
```

```

15 # The top two principal components
16 PC1_svd = X_scaled @ Vt.T[:, 0]
17 PC2_svd = X_scaled @ Vt.T[:, 1]
18
19 # Combine into a 2D array
20 X_pca_svd = np.vstack((PC1_svd, PC2_svd)).T
21
22 # Perform PCA using sklearn
23 pca = PCA(n_components=2)
24 X_pca_sklearn = pca.fit_transform(X_scaled)
25
26 # Check if the components are the same, up to a sign flip
27 sign_flip = np.sign(np.corrcoef(X_pca_svd[:, 0], X_pca_sklearn[:, 0])[0, 1])
28 X_pca_svd *= sign_flip
29
30 sign_flip = np.sign(np.corrcoef(X_pca_svd[:, 1], X_pca_sklearn[:, 1])[0, 1])
31 X_pca_svd[:, 1] *= sign_flip
32
33 assert np.allclose(X_pca_svd, X_pca_sklearn), "The two methods do not match!"
34
35 # Plotting
36 fig, (ax1, ax2) = plt.subplots(1, 2, figsize=(12, 6))
37
38 # SVD-based PCA scatter plot
39 ax1.scatter(X_pca_svd[y == 0, 0], X_pca_svd[y == 0, 1], label='Malignant', alpha=0.5)
40 ax1.scatter(X_pca_svd[y == 1, 0], X_pca_svd[y == 1, 1], label='Benign', alpha=0.5)
41 ax1.set_title('SVD-based PCA')
42 ax1.set_xlabel('First Principal Component')
43 ax1.set_ylabel('Second Principal Component')
44 ax1.legend()
45
46 # sklearn-based PCA scatter plot
47 ax2.scatter(X_pca_sklearn[y == 0, 0], X_pca_sklearn[y == 0, 1], label='Malignant', alpha=0.5)
48 ax2.scatter(X_pca_sklearn[y == 1, 0], X_pca_sklearn[y == 1, 1], label='Benign', alpha=0.5)
49 ax2.set_title('sklearn PCA')
50 ax2.set_xlabel('First Principal Component')
51 ax2.set_ylabel('Second Principal Component')
52 ax2.legend()
53
54 plt.show()

```

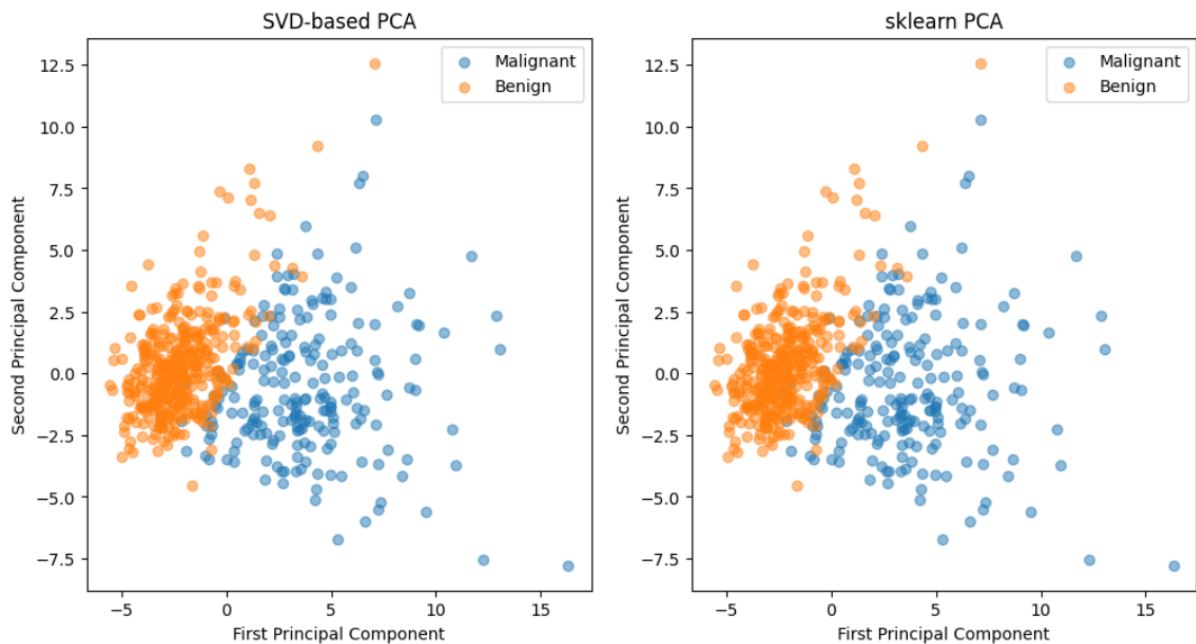


Figure 3: PCA/SVT Comparison