

```
In [1]: import matplotlib.pyplot as plt
import numpy as np
from numpy.linalg import svd
from sklearn import datasets
```

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 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.

```
In [2]: data = datasets.fetch_olivetti_faces()
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 ℓ_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 .

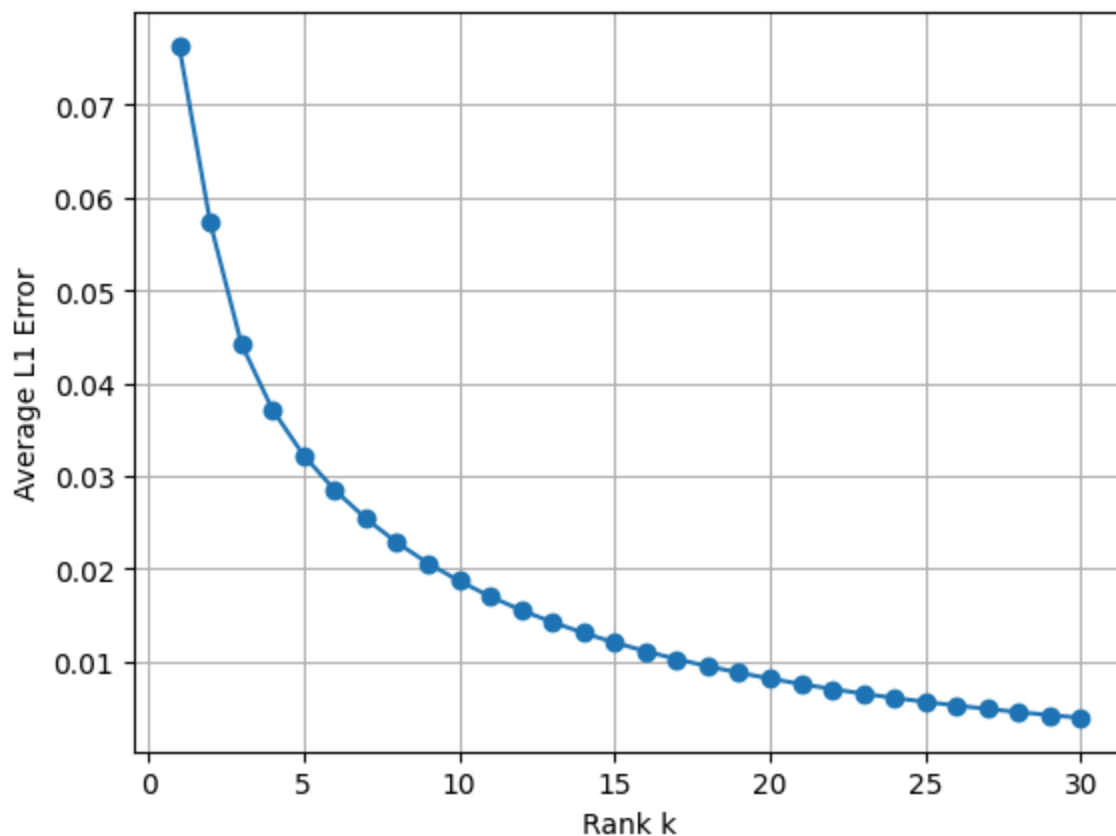
```
In [3]: # Returns the best rank-k approximation to M
def svd_reconstruct(M, k):
    U, S, VT = svd(M, full_matrices=False)
    Uk = U[:, :k]
    Sk = np.diag(S[:k])
    VTk = VT[:k, :]
    return Uk @ Sk @ VTk

def l1_error(M, M_approx):
    return np.mean(np.abs(M - M_approx))

errs = []
for k in range(1, 31):
    err = 0
    for img in images:
        reconstructed = svd_reconstruct(img, k)
        err += l1_error(img, reconstructed)
    errs.append(err/len(images))
```

```
In [4]: plt.plot(range(1, 31), errs, marker='o')
plt.xlabel('Rank k')
plt.ylabel('Average L1 Error')
```

```
plt.grid(True)
plt.show()
```



(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.

```
In [5]: np.random.seed(2025) # For reproducibility

random_img = images[np.random.randint(0, len(images))]

k = [10, 20, 30, 40]
fig, axes = plt.subplots(1, len(k) + 1, figsize=(15, 5))

axes[0].imshow(random_img, cmap='gray')
axes[0].set_title('Original Image')
axes[0].axis('off')

for i, rank in enumerate(k):
    axes[i + 1].imshow(svd_reconstruct(random_img, rank), cmap='gray')
    axes[i + 1].set_title(f'Rank {rank} Approx')
    axes[i + 1].axis('off')
```



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 mean 0 and standard deviation 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?

```
In [6]: cancer = datasets.load_breast_cancer()
```

```
In [7]: from sklearn.preprocessing import StandardScaler
        from sklearn.decomposition import PCA
```

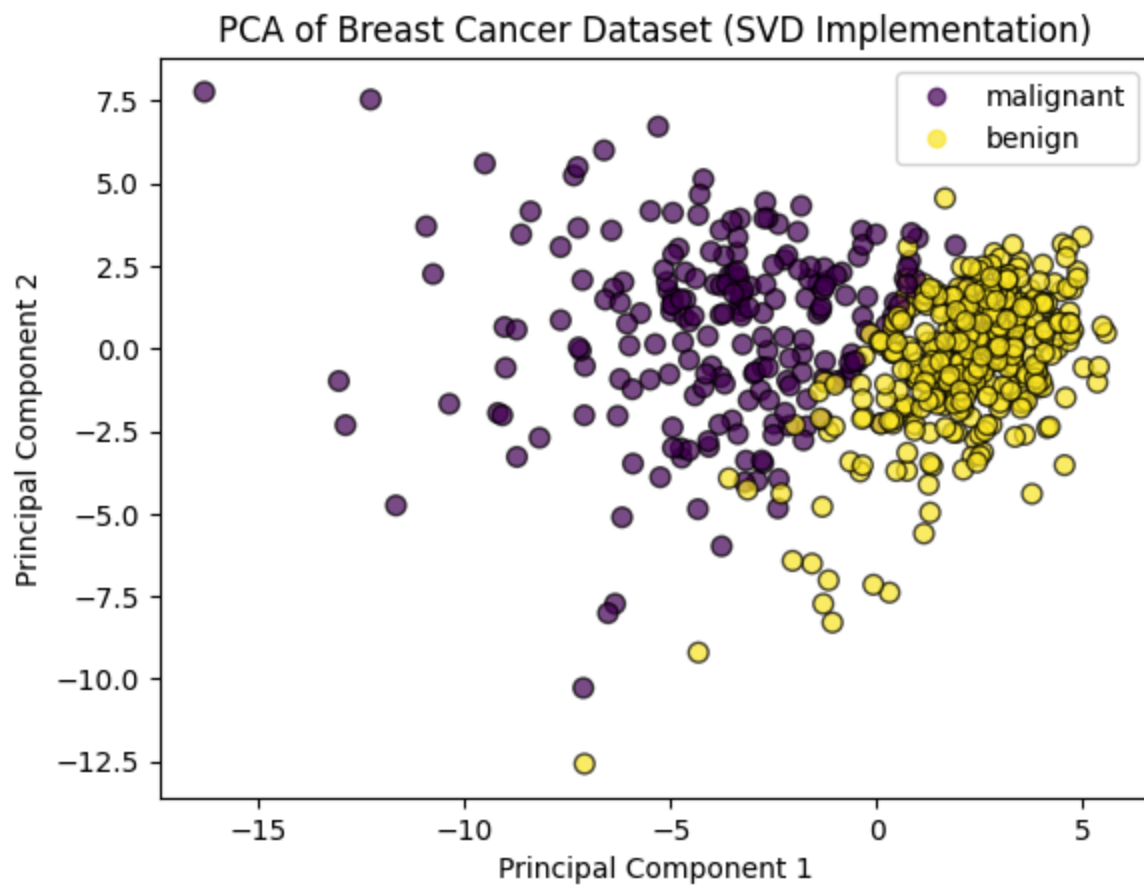
```
In [8]: x = cancer.data
        y = cancer.target
        target_names = cancer.target_names

        scaler = StandardScaler()
        x_scaled = scaler.fit_transform(x)
```

```
In [9]: #SVD Implementation (Not using PCA from sklearn)
        U, S, VT = svd(x_scaled, full_matrices=False)
        principal_components = VT.T[:, :2]
        x_pca_svd = x_scaled @ principal_components

        # Scatter plot of the two principal components
        sct = plt.scatter(x_pca_svd[:, 0], x_pca_svd[:, 1], c=y, edgecolors='k', s=50, cmap='vir
        plt.xlabel('Principal Component 1')
        plt.ylabel('Principal Component 2')
        plt.title('PCA of Breast Cancer Dataset (SVD Implementation)')
        plt.legend(handles=sct.legend_elements()[0], labels=list(target_names))
```

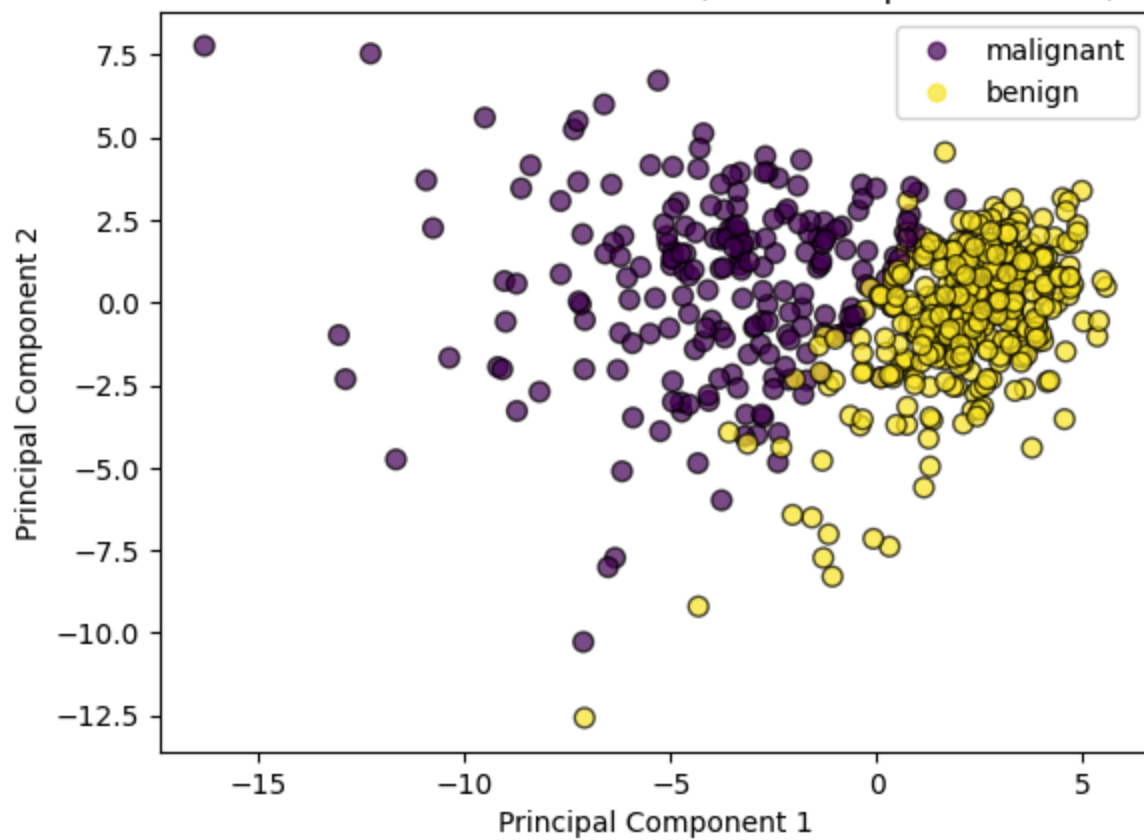
```
Out[9]: <matplotlib.legend.Legend at 0x127496b10>
```



```
In [10]: # PCA approach using sklearn
pca = PCA(n_components=2)
x_pca = pca.fit_transform(x_scaled)

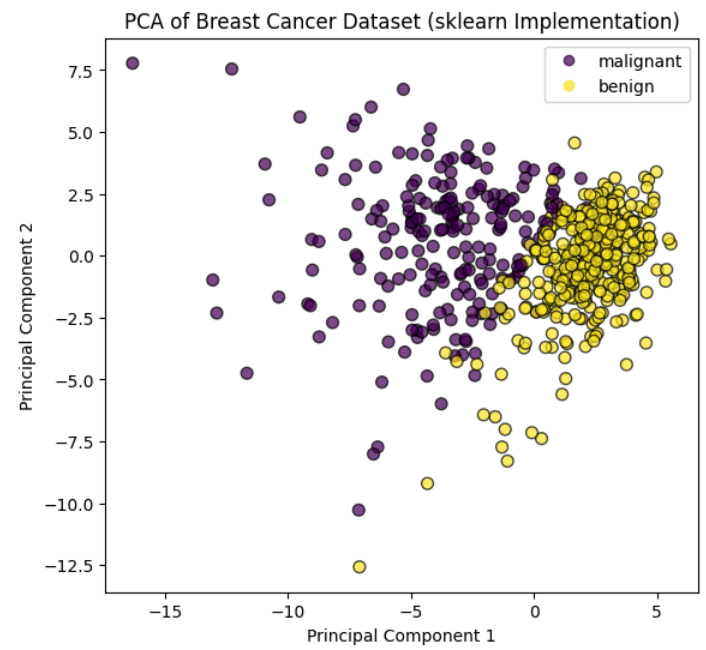
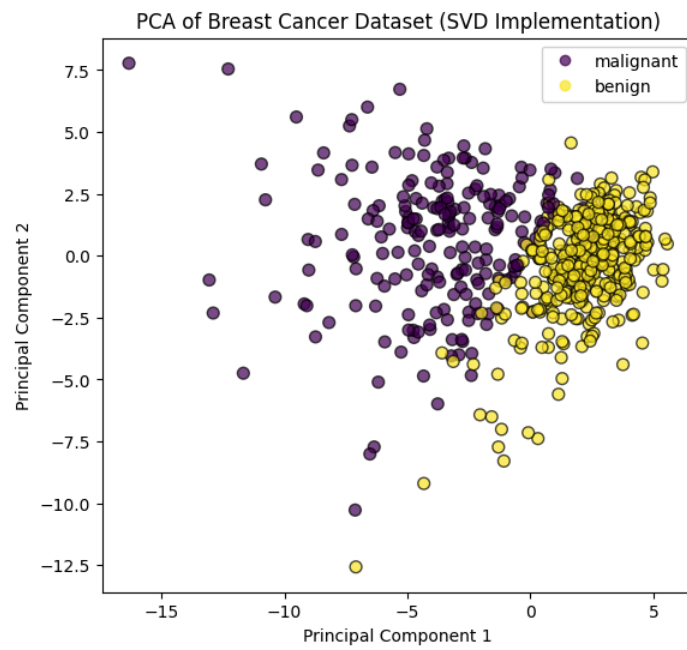
# Scatter plot of the two principal components
sct = plt.scatter(-x_pca[:, 0], -x_pca[:, 1], c=y, edgecolors='k', s=50, cmap='viridis',
plt.xlabel('Principal Component 1')
plt.ylabel('Principal Component 2')
plt.title('PCA of Breast Cancer Dataset (sklearn Implementation)')
plt.legend(handles=sct.legend_elements()[0], labels=list(target_names))
plt.show()
```

PCA of Breast Cancer Dataset (sklearn Implementation)



```
In [11]: # side by side comparison
fig, axes = plt.subplots(1, 2, figsize=(14, 6))
# SVD Implementation
sct1 = axes[0].scatter(x_pca_svd[:, 0], x_pca_svd[:, 1], c=y, edgecolors='k', s=50, cmap=
axes[0].set_xlabel('Principal Component 1')
axes[0].set_ylabel('Principal Component 2')
axes[0].set_title('PCA of Breast Cancer Dataset (SVD Implementation)')
axes[0].legend(handles=sct1.legend_elements()[0], labels=list(target_names))

# sklearn PCA Implementation
sct2 = axes[1].scatter(-x_pca[:, 0], -x_pca[:, 1], c=y, edgecolors='k', s=50, cmap='viri
axes[1].set_xlabel('Principal Component 1')
axes[1].set_ylabel('Principal Component 2')
axes[1].set_title('PCA of Breast Cancer Dataset (sklearn Implementation)')
axes[1].legend(handles=sct2.legend_elements()[0], labels=list(target_names))
plt.show()
```



The data does appear to be separable in 2 dimensions