# eda-lab-6-b

August 19, 2024

```python
[30]: import numpy as np
      import pandas as pd
```

```python
[31]: data = pd.read_csv('Data.csv')
      data.head()
```

```
[31]:          id diagnosis  radius_mean  texture_mean  perimeter_mean  area_mean  \
      0    842302         M        17.99         10.38          122.80     1001.0
      1    842517         M        20.57         17.77          132.90     1326.0
      2  84300903         M        19.69         21.25          130.00     1203.0
      3  84348301         M        11.42         20.38           77.58      386.1
      4  84358402         M        20.29         14.34          135.10     1297.0

         smoothness_mean  compactness_mean  concavity_mean  concave points_mean  \
      0          0.11840           0.27760          0.3001              0.14710
      1          0.08474           0.07864          0.0869              0.07017
      2          0.10960           0.15990          0.1974              0.12790
      3          0.14250           0.28390          0.2414              0.10520
      4          0.10030           0.13280          0.1980              0.10430

         …  texture_worst  perimeter_worst  area_worst  smoothness_worst  \
      0  …          17.33           184.60      2019.0            0.1622
      1  …          23.41           158.80      1956.0            0.1238
      2  …          25.53           152.50      1709.0            0.1444
      3  …          26.50            98.87       567.7            0.2098
      4  …          16.67           152.20      1575.0            0.1374

         compactness_worst  concavity_worst  concave points_worst  symmetry_worst  \
      0             0.6656           0.7119                0.2654          0.4601
      1             0.1866           0.2416                0.1860          0.2750
      2             0.4245           0.4504                0.2430          0.3613
      3             0.8663           0.6869                0.2575          0.6638
      4             0.2050           0.4000                0.1625          0.2364

         fractal_dimension_worst  Unnamed: 32
      0                  0.11890          NaN
      1                  0.08902          NaN
```

```
2                0.08758        NaN
3                0.17300        NaN
4                0.07678        NaN

[5 rows x 33 columns]
```

[32]: 
```python
data = data.drop(columns='Unnamed: 32')
```

[33]: 
```python
X = data.drop(columns='diagnosis')
y = [0 if row == 'B' else 1 for row in data['diagnosis']]
```
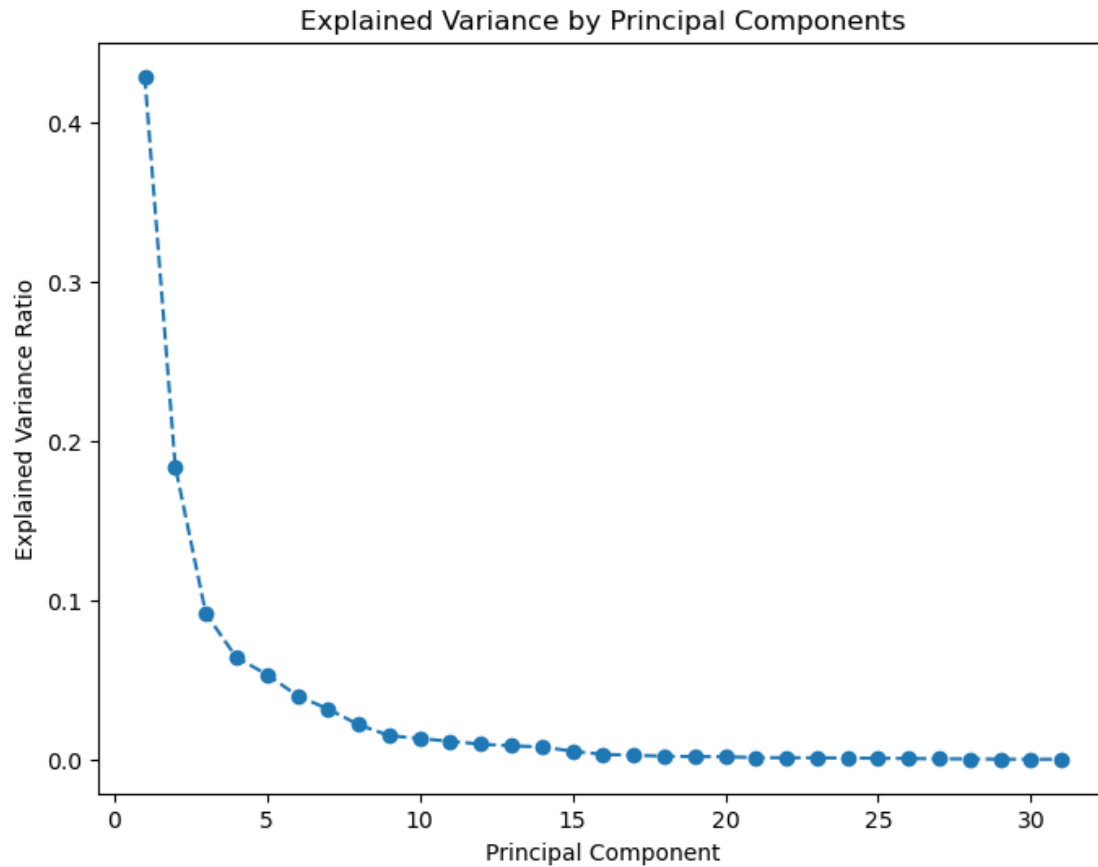
[34]: 
```python
from sklearn.preprocessing import StandardScaler
scaler = StandardScaler()
X_scaled = scaler.fit_transform(X)
```

[35]: 
```python
from sklearn.decomposition import PCA
import matplotlib.pyplot as plt

pca = PCA()
X_pca = pca.fit_transform(X_scaled)

explained_variance = pca.explained_variance_ratio_

# Plot the explained variance
plt.figure(figsize=(8, 6))
plt.plot(range(1, len(explained_variance) + 1), explained_variance, marker='o',
  linestyle='--')
plt.title('Explained Variance by Principal Components')
plt.xlabel('Principal Component')
plt.ylabel('Explained Variance Ratio')
plt.show()
```

## Explained Variance by Principal Components



[37]:
```python
pca_2 = PCA(n_components=2)
X_pca_2 = pca_2.fit_transform(X_scaled)

plt.figure(figsize=(8, 6))
plt.scatter(X_pca_2[:, 0], X_pca_2[:, 1], c=y, cmap='viridis', edgecolor='k',␣
 ↪s=50)
plt.title('PCA Scatterplot (First 2 Principal Components)')
plt.xlabel('Principal Component 1')
plt.ylabel('Principal Component 2')
plt.colorbar(label='Diagnosis (0 = Benign, 1 = Malignant)')
plt.show()
```

PCA Scatterplot (First 2 Principal Components)