

▼ Description

This Python notebook aims to provide a hands-on experience with Principal Component Analysis (PCA) and Linear Discriminant Analysis (LDA) using NumPy. It focuses on allowing for a clear understanding of their mechanics and resulting data transformations.

- PCA with numpy
- LDA with numpy

```
import numpy as np
import pandas as pd
from matplotlib import pyplot as plt
import seaborn as sns
```

```
df = pd.read_csv('https://archive.ics.uci.edu/ml/'
                 'machine-learning-databases/wine/wine.data',
                 header=None)

df.columns = ['Class label', 'Alcohol', 'Malic acid', 'Ash',
              'Alcalinity of ash', 'Magnesium', 'Total phenols',
              'Flavanoids', 'Nonflavanoid phenols', 'Proanthocyanins',
              'Color intensity', 'Hue',
              'OD280/OD315 of diluted wines', 'Proline']
```

```
df.head()
```



	Class label	Alcohol	Malic acid	Ash	Alcalinity of ash	Magnesium	Total phenols	Flavanoids	Nonflavanoid phenols	Proanthocyanins	Color intensity	Hue	OD280/OD315 of diluted wines	Proline
0	1	14.23	1.71	2.43	15.6	127	2.80	3.06	0.28	2.29	5.64	1.04	3.92	
1	1	13.20	1.78	2.14	11.2	100	2.65	2.76	0.26	1.28	4.38	1.05	3.40	
2	1	13.16	2.36	2.67	18.6	101	2.80	3.24	0.30	2.81	5.68	1.03	3.17	
3	1	14.37	1.95	2.50	16.8	113	3.85	3.49	0.24	2.18	7.80	0.86	3.45	
4	1	13.24	2.59	2.87	21.0	118	2.80	2.69	0.39	1.82	4.32	1.04	2.93	

Next steps:

[Generate code with df](#)

[View recommended plots](#)

[New interactive sheet](#)

```
df.describe()
```



	Class label	Alcohol	Malic acid	Ash	Alcalinity of ash	Magnesium	Total phenols	Flavanoids	Nonflavanoid phenols	Proanthocyanins	Color intensity	Hue	OD280/OD315 of diluted wines	Proline
count	178.000000	178.000000	178.000000	178.000000	178.000000	178.000000	178.000000	178.000000	178.000000	178.000000	178.000000	178.000000	178.000000	178.000000
mean	1.938202	13.000618	2.336348	2.366517	19.494944	99.741573	2.295112	2.029270	0.361854	1.590899	5.012548	1.034470	3.393411	1596.019843
std	0.775035	0.811827	1.117146	0.274344	3.339564	14.282484	0.625851	0.998859	0.124453	0.572359	2.310384	0.185699	0.705553	752.214749
min	1.000000	11.030000	0.740000	1.360000	10.600000	70.000000	0.980000	0.340000	0.130000	0.410000	1.250000	0.830000	2.930000	417.000000
25%	1.000000	12.362500	1.602500	2.210000	17.200000	88.000000	1.742500	1.205000	0.270000	1.250000	3.393411	0.930000	3.170000	1195.000000
50%	2.000000	13.050000	1.865000	2.360000	19.500000	98.000000	2.355000	2.135000	0.340000	1.555000	4.612548	1.030000	3.393411	1596.019843
75%	3.000000	13.677500	3.082500	2.557500	21.500000	107.000000	2.800000	2.875000	0.437500	1.950000	6.250000	1.040000	3.920000	1707.000000
max	3.000000	14.830000	5.800000	3.230000	30.000000	162.000000	3.880000	5.080000	0.660000	3.580000	13.000000	1.050000	3.920000	1707.000000

```
df.info()
```



```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 178 entries, 0 to 177
Data columns (total 14 columns):
#   Column              Non-Null Count  Dtype
---  -
0   Class label         178 non-null    int64
1   Alcohol             178 non-null    float64
```

```

2 Malic acid 178 non-null float64
3 Ash 178 non-null float64
4 Alcalinity of ash 178 non-null float64
5 Magnesium 178 non-null int64
6 Total phenols 178 non-null float64
7 Flavanoids 178 non-null float64
8 Nonflavanoid phenols 178 non-null float64
9 Proanthocyanins 178 non-null float64
10 Color intensity 178 non-null float64
11 Hue 178 non-null float64
12 OD280/OD315 of diluted wines 178 non-null float64
13 Proline 178 non-null int64
dtypes: float64(11), int64(3)
memory usage: 19.6 KB

```

```

from sklearn.model_selection import train_test_split

X, y = df.iloc[:, 1:].values, df.iloc[:, 0].values

X_train, X_test, y_train, y_test = \
    train_test_split(X, y, test_size=0.3,
                    stratify=y,
                    random_state=0)

```

```

from sklearn.preprocessing import StandardScaler

sc = StandardScaler()
X_train_std = sc.fit_transform(X_train)
X_test_std = sc.transform(X_test)

```

✓ 1.PCA with numpy

```

# Covariance matrix decomposition into eigenval, eigenvect
import numpy as np
cov_mat = np.cov(X_train_std.T)
eigen_vals, eigen_vecs = np.linalg.eig(cov_mat)

print('\nEigenvalues \n', eigen_vals)

```



```

Eigenvalues
[4.84274532 2.41602459 1.54845825 0.96120438 0.84166161 0.6620634
 0.51828472 0.34650377 0.3131368 0.10754642 0.21357215 0.15362835
 0.1808613 ]

```

✓ Eigenvalue (Explained Variance)

```

# Total and explained variance
tot = sum(eigen_vals)
var_exp = [(i / tot) for i in sorted(eigen_vals, reverse=True)]
cum_var_exp = np.cumsum(var_exp)

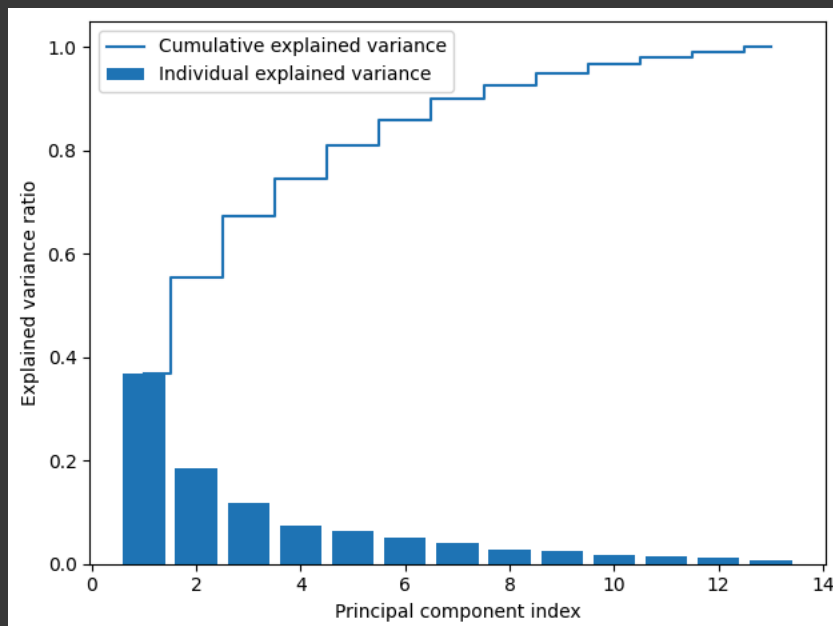
```

```

plt.bar(range(1, 14), var_exp, align='center',
        label='Individual explained variance')
plt.step(range(1, 14), cum_var_exp, where='mid',
        label='Cumulative explained variance')
plt.ylabel('Explained variance ratio')
plt.xlabel('Principal component index')
plt.legend(loc='best')
plt.tight_layout()

plt.show()

```



✓ Eigenvector (Principal Components)

```
# Make a list of (eigenvalue, eigenvector) tuples
eigen_pairs = [(np.abs(eigen_vals[i]), eigen_vecs[:, i])
               for i in range(len(eigen_vals))]
```

```
# Sort the (eigenvalue, eigenvector) tuples from high to low
eigen_pairs.sort(key=lambda k: k[0], reverse=True)
```

```
w = np.hstack((eigen_pairs[0][1][:, np.newaxis],
               eigen_pairs[1][1][:, np.newaxis]))
print('Matrix W:\n', w)
```



```
Matrix W:
[[-0.13724218  0.50303478]
 [ 0.24724326  0.16487119]
 [-0.02545159  0.24456476]
 [ 0.20694508 -0.11352904]
 [-0.15436582  0.28974518]
 [-0.39376952  0.05080104]
 [-0.41735106 -0.02287338]
 [ 0.30572896  0.09048885]
 [-0.30668347  0.00835233]
 [ 0.07554066  0.54977581]
 [-0.32613263 -0.20716433]
 [-0.36861022 -0.24902536]
 [-0.29669651  0.38022942]]
```

```
# how the first data point converted
# xW = x'
X_train_std[0].dot(w)
```



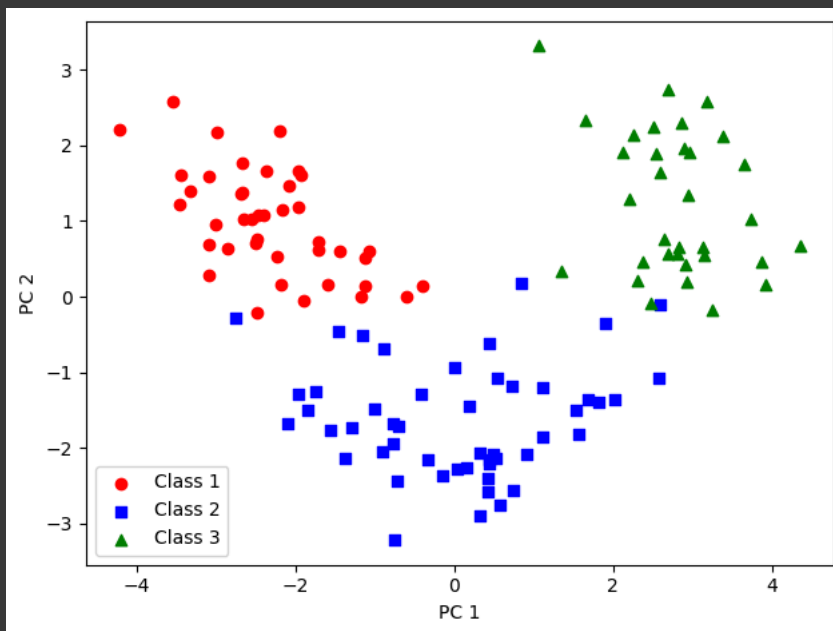
```
array([2.38299011, 0.45458499])
```

```
X_train_pca = X_train_std.dot(w)
colors = ['r', 'b', 'g']
markers = ['o', 's', '^']

for l, c, m in zip(np.unique(y_train), colors, markers):
    plt.scatter(X_train_pca[y_train == l, 0],
                X_train_pca[y_train == l, 1],
                c=c, label=f'Class {l}', marker=m)

plt.xlabel('PC 1')
plt.ylabel('PC 2')
plt.legend(loc='lower left')
plt.tight_layout()
```

```
plt.show()
```



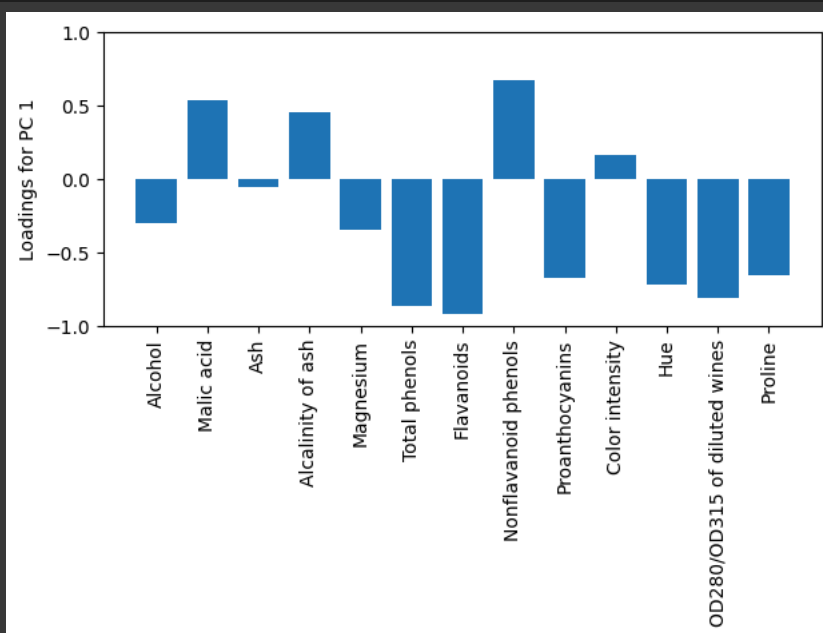
```
# Assess the feature contribution on first principal component
loadings = eigen_vecs * np.sqrt(eigen_vals)
```

```
fig, ax = plt.subplots()
```

```
ax.bar(range(13), loadings[:, 0], align='center')
ax.set_ylabel('Loadings for PC 1')
ax.set_xticks(range(13))
ax.set_xticklabels(df.columns[1:], rotation=90)
```

```
plt.ylim([-1, 1])
plt.tight_layout()
```

```
plt.show()
```



```
loadings[:, 0]
```



```
array([-0.3020184 ,  0.54408942, -0.05600938,  0.45540829, -0.33970111,
        -0.8665386 , -0.9184327 ,  0.67279444, -0.67489496,  0.16623657,
        -0.71769524, -0.81117245, -0.65291742])
```

✓ 2.LDA with numpy

```
# Calculate the feature means by class
np.set_printoptions(precision=4)
```

```
mean_vecs = []
for label in range(1, 4):
    mean_vecs.append(np.mean(X_train_std[y_train == label], axis=0))
    print(f'MV {label}: {mean_vecs[label - 1]}\n')
```

➦ MV 1: [0.9066 -0.3497 0.3201 -0.7189 0.5056 0.8807 0.9589 -0.5516 0.5416
0.2338 0.5897 0.6563 1.2075]

MV 2: [-0.8749 -0.2848 -0.3735 0.3157 -0.3848 -0.0433 0.0635 -0.0946 0.0703
-0.8286 0.3144 0.3608 -0.7253]

MV 3: [0.1992 0.866 0.1682 0.4148 -0.0451 -1.0286 -1.2876 0.8287 -0.7795
0.9649 -1.209 -1.3622 -0.4013]

```
# Compute the within-class scatter matrix Sw
d = 13 # number of features
S_W = np.zeros((d, d))
for label, mv in zip(range(1, 4), mean_vecs):
    class_scatter = np.zeros((d, d)) # scatter matrix for each class
    for row in X_train_std[y_train == label]:
        row, mv = row.reshape(d, 1), mv.reshape(d, 1) # make column vectors
        class_scatter += (row - mv).dot((row - mv).T)
    S_W += class_scatter # sum class scatter matrices

print('Within-class scatter matrix: '
      f'{S_W.shape[0]}x{S_W.shape[1]}')
```

➦ Within-class scatter matrix: 13x13

```
# Check number of samples for each class
print('Class label distribution:',
      np.bincount(y_train)[1:])
```

➦ Class label distribution: [41 50 33]

```
# Scale within-class scatter matrix
d = 13 # number of features
S_W = np.zeros((d, d))
for label, mv in zip(range(1, 4), mean_vecs):
    class_scatter = np.cov(X_train_std[y_train == label].T)
    S_W += class_scatter

print('Scaled within-class scatter matrix: '
      f'{S_W.shape[0]}x{S_W.shape[1]}')
```

➦ Scaled within-class scatter matrix: 13x13

```
# Compute between-class scatter matrix
mean_overall = np.mean(X_train_std, axis=0)
mean_overall = mean_overall.reshape(d, 1) # make column vector

d = 13 # number of features
S_B = np.zeros((d, d))

for i, mean_vec in enumerate(mean_vecs):
    n = X_train_std[y_train == i + 1, :].shape[0]
    mean_vec = mean_vec.reshape(d, 1) # make column vector
    S_B += n * (mean_vec - mean_overall).dot((mean_vec - mean_overall).T)

print('Between-class scatter matrix: '
      f'{S_B.shape[0]}x{S_B.shape[1]}')
```

➦ Between-class scatter matrix: 13x13

```
# Decompose the scatter matrix: Sw^-1*Sb to eigenval, eigenvec
eigen_vals, eigen_vecs = np.linalg.eig(np.linalg.inv(S_W).dot(S_B))
```

```
# Make a list of (eigenvalue, eigenvector) tuples
eigen_pairs = [(np.abs(eigen_vals[i]), eigen_vecs[:, i])
               for i in range(len(eigen_vals))]

# Sort the (eigenvalue, eigenvector) tuples from high to low
eigen_pairs = sorted(eigen_pairs, key=lambda k: k[0], reverse=True)

# Visually confirm that the list is correctly sorted by decreasing eigenvalues

print('Eigenvalues in descending order:\n')
for eigen_val in eigen_pairs:
    print(eigen_val[0])
```

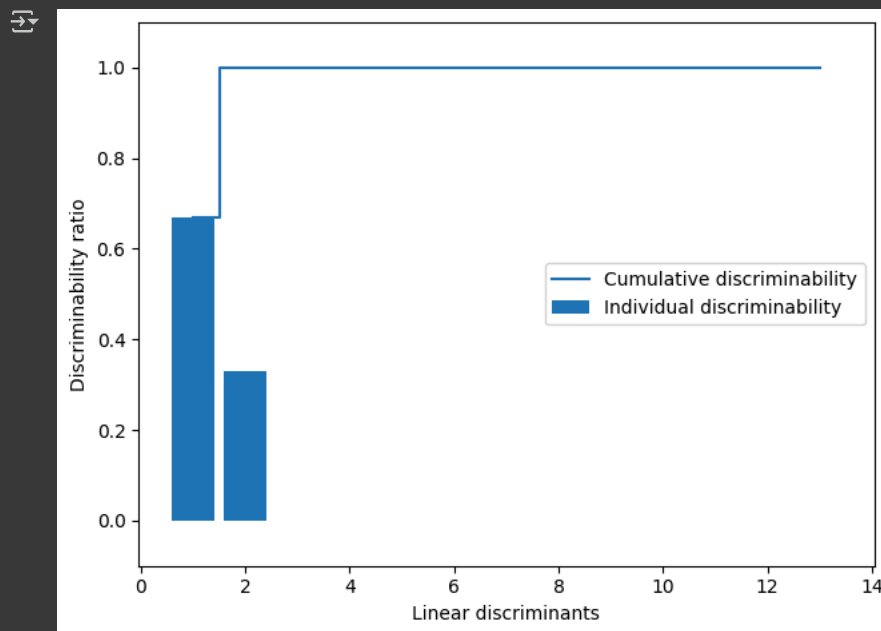
Eigenvalues in descending order:

```
349.61780890599397
172.7615221897938
3.342838214841367e-14
2.842170943040401e-14
2.5545786180111422e-14
1.7533939180734234e-14
1.7533939180734234e-14
1.6579193995960903e-14
1.6579193995960903e-14
8.242524002707225e-15
8.242524002707225e-15
6.36835506006027e-15
2.974634375545734e-15
```

```
tot = sum(eigen_vals.real)
discr = [(i / tot) for i in sorted(eigen_vals.real, reverse=True)]
cum_discr = np.cumsum(discr)

plt.bar(range(1, 14), discr, align='center',
        label='Individual discriminability')
plt.step(range(1, 14), cum_discr, where='mid',
        label='Cumulative discriminability')
plt.ylabel('Discriminability ratio')
plt.xlabel('Linear discriminants')
plt.ylim([-0.1, 1.1])
plt.legend(loc='best')
plt.tight_layout()

plt.show()
```



```
w = np.hstack((eigen_pairs[0][1][:, np.newaxis].real,
               eigen_pairs[1][1][:, np.newaxis].real))
print('Matrix W:\n', w)
```

Matrix W:

```
[[-0.1481 -0.4092]
 [ 0.0908 -0.1577]]
```

```

[-0.0168 -0.3537]
[ 0.1484  0.3223]
[-0.0163 -0.0817]
[ 0.1913  0.0842]
[-0.7338  0.2823]
[-0.075  -0.0102]
[ 0.0018  0.0907]
[ 0.294  -0.2152]
[-0.0328  0.2747]
[-0.3547 -0.0124]
[-0.3915 -0.5958]]

```

```

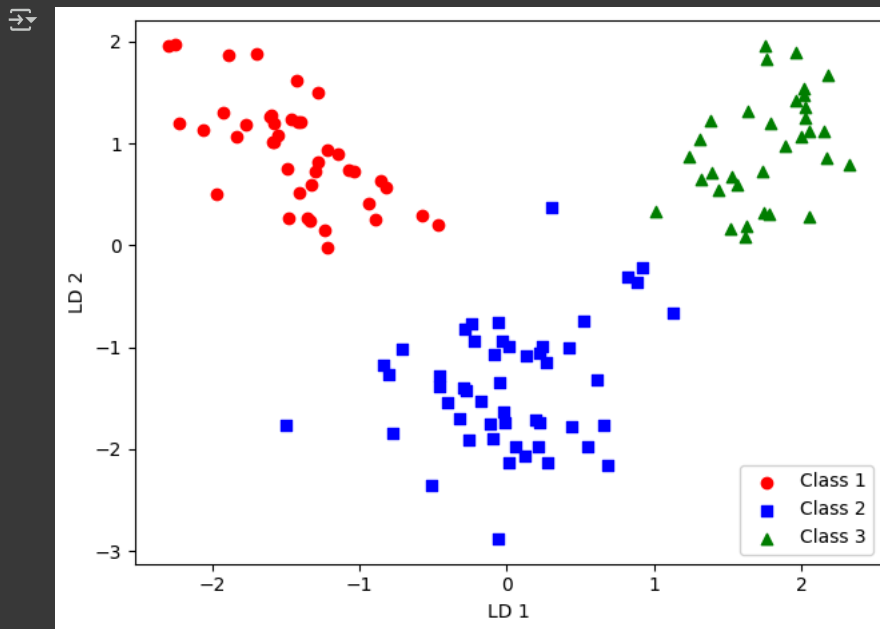
# Project the data points to new feature space
X_train_lda = X_train_std.dot(w)
colors = ['r', 'b', 'g']
markers = ['o', 's', '^']

for l, c, m in zip(np.unique(y_train), colors, markers):
    plt.scatter(X_train_lda[y_train == l, 0],
                X_train_lda[y_train == l, 1] * (-1),
                c=c, label=f'Class {l}', marker=m)

plt.xlabel('LD 1')
plt.ylabel('LD 2')
plt.legend(loc='lower right')
plt.tight_layout()

plt.show()

```



✓ Your work

Objective: learn how to apply **sklearn** PCA and LDA APIs

Task: Please review the sklearn code in [this link](#) and apply it to [Diabetes dataset](#). Complete the data preprocessing(e.g. features standardization) prior to dimension reduction.

For an easier start, you can also reference to the official example on sklearn.

- [PCA](#)
- [LDA](#)

Please submit your work to BrightSpace by 3/16 11:59 pm.

```

import numpy as np
import pandas as pd
from sklearn.model_selection import train_test_split
from sklearn.preprocessing import StandardScaler
from sklearn.decomposition import PCA
from sklearn.discriminant_analysis import LinearDiscriminantAnalysis as LDA

```

```
import matplotlib.pyplot as plt
import seaborn as sns
```

```
url = 'https://raw.githubusercontent.com/mdogy/dataForEng1999/master/pi_diabetes.csv'
df = pd.read_csv(url)
```

```
print(df.head())
print(df.describe())
print(df.info())
```

```

Pregnancies  Glucose  BloodPressure  SkinThickness  Insulin  BMI  \
0            6       148             72           35      0  33.6
1            1        85             66           29      0  26.6
2            8       183             64            0      0  23.3
3            1        89             66           23     94  28.1
4            0       137             40           35    168  43.1
```

```

DiabetesPedigreeFunction  Age  Outcome
0            0.627      50         1
1            0.351      31         0
2            0.672      32         1
3            0.167      21         0
4            2.288      33         1
```

```

Pregnancies  Glucose  BloodPressure  SkinThickness  Insulin  \
count  768.000000  768.000000  768.000000  768.000000  768.000000
mean    3.845052  120.894531   69.105469   20.536458   79.799479
std     3.369578   31.972618   19.355807   15.952218  115.244002
min     0.000000   0.000000    0.000000    0.000000    0.000000
25%     1.000000   99.000000   62.000000    0.000000    0.000000
50%     3.000000  117.000000   72.000000   23.000000   30.500000
75%     6.000000  140.250000   80.000000   32.000000  127.250000
max    17.000000  199.000000  122.000000   99.000000  846.000000
```

```

BMI  DiabetesPedigreeFunction  Age  Outcome
count  768.000000              768.000000  768.000000  768.000000
mean    31.992578              0.471876   33.240885    0.348958
std     7.884160              0.331329   11.760232    0.476951
min     0.000000              0.078000   21.000000    0.000000
25%    27.300000              0.243750   24.000000    0.000000
50%    32.000000              0.372500   29.000000    0.000000
75%    36.600000              0.626250   41.000000    1.000000
max    67.100000              2.420000   81.000000    1.000000
```

```
<class 'pandas.core.frame.DataFrame'>
```

```
RangeIndex: 768 entries, 0 to 767
```

```
Data columns (total 9 columns):
```

```

#  Column  Non-Null Count  Dtype
---  -
0  Pregnancies  768 non-null  int64
1  Glucose      768 non-null  int64
2  BloodPressure  768 non-null  int64
3  SkinThickness  768 non-null  int64
4  Insulin      768 non-null  int64
5  BMI          768 non-null  float64
6  DiabetesPedigreeFunction  768 non-null  float64
7  Age          768 non-null  int64
8  Outcome      768 non-null  int64
```

```
dtypes: float64(2), int64(7)
```

```
memory usage: 54.1 KB
```

```
None
```

```
# Replace zero values with NaN
```

```
df[['Glucose', 'BloodPressure', 'SkinThickness', 'Insulin', 'BMI']] = df[['Glucose', 'BloodPressure', 'SkinThickness', 'Insulin', 'BMI']].re
```

```
# Impute missing values with the median of each column
```

```
df.fillna(df.median(), inplace=True)
```

```
X = df.drop('Outcome', axis=1).values
```

```
y = df['Outcome'].values
```

```
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3, stratify=y, random_state=0)
```

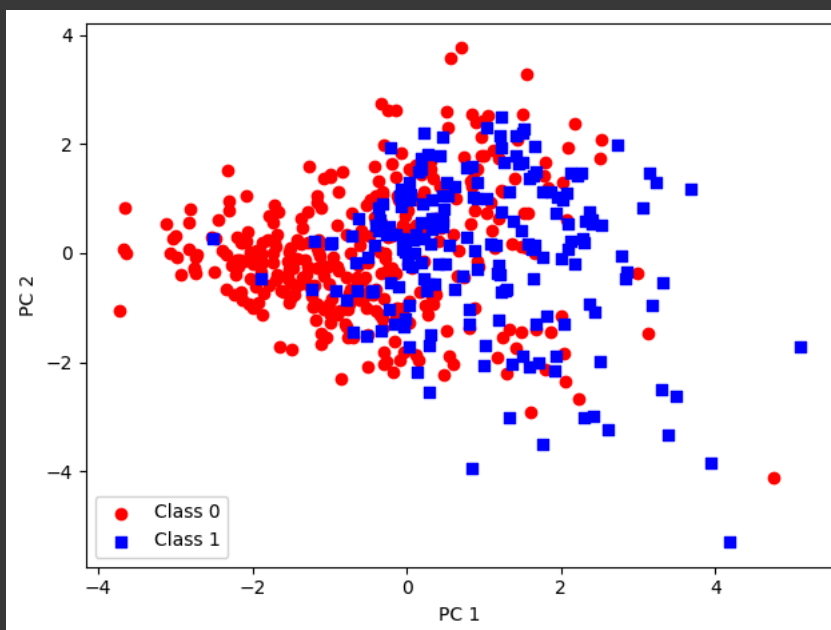


```
sc = StandardScaler()
X_train_std = sc.fit_transform(X_train)
X_test_std = sc.transform(X_test)
```

```
pca = PCA(n_components=2)
X_train_pca = pca.fit_transform(X_train_std)
X_test_pca = pca.transform(X_test_std)
```

```
colors = ['r', 'b']
markers = ['o', 's']

for l, c, m in zip(np.unique(y_train), colors, markers):
    plt.scatter(X_train_pca[y_train == l, 0], X_train_pca[y_train == l, 1], c=c, label=f'Class {l}', marker=m)
plt.xlabel('PC 1')
plt.ylabel('PC 2')
plt.legend(loc='lower left')
plt.tight_layout()
plt.show()
```



```
lda = LDA(n_components=1)
X_train_lda = lda.fit_transform(X_train_std, y_train)
X_test_lda = lda.transform(X_test_std)
```

```
for l, c, m in zip(np.unique(y_train), colors, markers):
    plt.scatter(X_train_lda[y_train == l], np.zeros_like(X_train_lda[y_train == l]),
                c=c, label=f'Class {l}', marker=m)

plt.xlabel('LD 1')
plt.yticks([])
plt.legend(loc='lower left')
plt.tight_layout()
plt.show()
```

