

PCA

February 20, 2024

1 Import necessary libraries

```
[ ]: import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
import warnings
warnings.filterwarnings('ignore')
```

2 Load Dataset and basic exploration

```
[ ]: df = pd.read_csv("datasets/breast-cancer.csv")
df.head()
```

```
[ ]:
      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

      ...  radius_worst  texture_worst  perimeter_worst  area_worst  \
0  ...         25.38         17.33         184.60         2019.0
1  ...         24.99         23.41         158.80         1956.0
2  ...         23.57         25.53         152.50         1709.0
3  ...         14.91         26.50          98.87          567.7
4  ...         22.54         16.67         152.20         1575.0

      smoothness_worst  compactness_worst  concavity_worst  concave points_worst  \
0          0.1622         0.6656         0.7119         0.2654
```

1	0.1238	0.1866	0.2416	0.1860
2	0.1444	0.4245	0.4504	0.2430
3	0.2098	0.8663	0.6869	0.2575
4	0.1374	0.2050	0.4000	0.1625

	symmetry_worst	fractal_dimension_worst
0	0.4601	0.11890
1	0.2750	0.08902
2	0.3613	0.08758
3	0.6638	0.17300
4	0.2364	0.07678

[5 rows x 32 columns]

```
[ ]: df.columns
```

```
[ ]: Index(['id', 'diagnosis', 'radius_mean', 'texture_mean', 'perimeter_mean',
          'area_mean', 'smoothness_mean', 'compactness_mean', 'concavity_mean',
          'concave points_mean', 'symmetry_mean', 'fractal_dimension_mean',
          'radius_se', 'texture_se', 'perimeter_se', 'area_se', 'smoothness_se',
          'compactness_se', 'concavity_se', 'concave points_se', 'symmetry_se',
          'fractal_dimension_se', 'radius_worst', 'texture_worst',
          'perimeter_worst', 'area_worst', 'smoothness_worst',
          'compactness_worst', 'concavity_worst', 'concave points_worst',
          'symmetry_worst', 'fractal_dimension_worst'],
          dtype='object')
```

```
[ ]: df = df.drop(columns=["id"])
df.head()
```

```
[ ]:  diagnosis  radius_mean  texture_mean  perimeter_mean  area_mean  \
0         M         17.99         10.38         122.80        1001.0
1         M         20.57         17.77         132.90        1326.0
2         M         19.69         21.25         130.00        1203.0
3         M         11.42         20.38          77.58         386.1
4         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

    symmetry_mean  ...  radius_worst  texture_worst  perimeter_worst  \
0         0.2419  ...         25.38         17.33         184.60
1         0.1812  ...         24.99         23.41         158.80
```

2	0.2069	...	23.57	25.53	152.50
3	0.2597	...	14.91	26.50	98.87
4	0.1809	...	22.54	16.67	152.20

	area_worst	smoothness_worst	compactness_worst	concavity_worst	\
0	2019.0	0.1622	0.6656	0.7119	
1	1956.0	0.1238	0.1866	0.2416	
2	1709.0	0.1444	0.4245	0.4504	
3	567.7	0.2098	0.8663	0.6869	
4	1575.0	0.1374	0.2050	0.4000	

	concave points_worst	symmetry_worst	fractal_dimension_worst
0	0.2654	0.4601	0.11890
1	0.1860	0.2750	0.08902
2	0.2430	0.3613	0.08758
3	0.2575	0.6638	0.17300
4	0.1625	0.2364	0.07678

[5 rows x 31 columns]

```
[ ]: df.describe()
```

```
[ ]:
```

	radius_mean	texture_mean	perimeter_mean	area_mean	\
count	569.000000	569.000000	569.000000	569.000000	
mean	14.127292	19.289649	91.969033	654.889104	
std	3.524049	4.301036	24.298981	351.914129	
min	6.981000	9.710000	43.790000	143.500000	
25%	11.700000	16.170000	75.170000	420.300000	
50%	13.370000	18.840000	86.240000	551.100000	
75%	15.780000	21.800000	104.100000	782.700000	
max	28.110000	39.280000	188.500000	2501.000000	

	smoothness_mean	compactness_mean	concavity_mean	concave points_mean	\
count	569.000000	569.000000	569.000000	569.000000	
mean	0.096360	0.104341	0.088799	0.048919	
std	0.014064	0.052813	0.079720	0.038803	
min	0.052630	0.019380	0.000000	0.000000	
25%	0.086370	0.064920	0.029560	0.020310	
50%	0.095870	0.092630	0.061540	0.033500	
75%	0.105300	0.130400	0.130700	0.074000	
max	0.163400	0.345400	0.426800	0.201200	

	symmetry_mean	fractal_dimension_mean	...	radius_worst	\
count	569.000000	569.000000	...	569.000000	
mean	0.181162	0.062798	...	16.269190	
std	0.027414	0.007060	...	4.833242	
min	0.106000	0.049960	...	7.930000	

25%	0.161900	0.057700	...	13.010000
50%	0.179200	0.061540	...	14.970000
75%	0.195700	0.066120	...	18.790000
max	0.304000	0.097440	...	36.040000

	texture_worst	perimeter_worst	area_worst	smoothness_worst	\
count	569.000000	569.000000	569.000000	569.000000	
mean	25.677223	107.261213	880.583128	0.132369	
std	6.146258	33.602542	569.356993	0.022832	
min	12.020000	50.410000	185.200000	0.071170	
25%	21.080000	84.110000	515.300000	0.116600	
50%	25.410000	97.660000	686.500000	0.131300	
75%	29.720000	125.400000	1084.000000	0.146000	
max	49.540000	251.200000	4254.000000	0.222600	

	compactness_worst	concavity_worst	concave	points_worst	\
count	569.000000	569.000000		569.000000	
mean	0.254265	0.272188		0.114606	
std	0.157336	0.208624		0.065732	
min	0.027290	0.000000		0.000000	
25%	0.147200	0.114500		0.064930	
50%	0.211900	0.226700		0.099930	
75%	0.339100	0.382900		0.161400	
max	1.058000	1.252000		0.291000	

	symmetry_worst	fractal_dimension_worst
count	569.000000	569.000000
mean	0.290076	0.083946
std	0.061867	0.018061
min	0.156500	0.055040
25%	0.250400	0.071460
50%	0.282200	0.080040
75%	0.317900	0.092080
max	0.663800	0.207500

[8 rows x 30 columns]

```
[ ]: df.info()
```

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

```
RangeIndex: 569 entries, 0 to 568
```

```
Data columns (total 31 columns):
```

#	Column	Non-Null Count	Dtype
---	-----	-----	-----
0	diagnosis	569 non-null	object
1	radius_mean	569 non-null	float64
2	texture_mean	569 non-null	float64
3	perimeter_mean	569 non-null	float64

4	area_mean	569	non-null	float64
5	smoothness_mean	569	non-null	float64
6	compactness_mean	569	non-null	float64
7	concavity_mean	569	non-null	float64
8	concave points_mean	569	non-null	float64
9	symmetry_mean	569	non-null	float64
10	fractal_dimension_mean	569	non-null	float64
11	radius_se	569	non-null	float64
12	texture_se	569	non-null	float64
13	perimeter_se	569	non-null	float64
14	area_se	569	non-null	float64
15	smoothness_se	569	non-null	float64
16	compactness_se	569	non-null	float64
17	concavity_se	569	non-null	float64
18	concave points_se	569	non-null	float64
19	symmetry_se	569	non-null	float64
20	fractal_dimension_se	569	non-null	float64
21	radius_worst	569	non-null	float64
22	texture_worst	569	non-null	float64
23	perimeter_worst	569	non-null	float64
24	area_worst	569	non-null	float64
25	smoothness_worst	569	non-null	float64
26	compactness_worst	569	non-null	float64
27	concavity_worst	569	non-null	float64
28	concave points_worst	569	non-null	float64
29	symmetry_worst	569	non-null	float64
30	fractal_dimension_worst	569	non-null	float64

dtypes: float64(30), object(1)
memory usage: 137.9+ KB

```
[ ]: df.isnull().sum()
```

```
[ ]: diagnosis      0
      radius_mean   0
      texture_mean   0
      perimeter_mean 0
      area_mean      0
      smoothness_mean 0
      compactness_mean 0
      concavity_mean  0
      concave points_mean 0
      symmetry_mean    0
      fractal_dimension_mean 0
      radius_se        0
      texture_se        0
      perimeter_se      0
      area_se           0
```

```

smoothness_se          0
compactness_se         0
concavity_se           0
concave points_se      0
symmetry_se            0
fractal_dimension_se   0
radius_worst           0
texture_worst          0
perimeter_worst        0
area_worst             0
smoothness_worst       0
compactness_worst      0
concavity_worst        0
concave points_worst   0
symmetry_worst         0
fractal_dimension_worst 0
dtype: int64

```

2.0.1 There are no null values from the above

3 Encoding the target column as it's object

```

[ ]: from sklearn.preprocessing import LabelEncoder
le = LabelEncoder()
df["diagnosis"] = le.fit_transform(df["diagnosis"])

```

4 Dependent and Independent Variable Split && Train Test Split

```

[ ]: x = df.iloc[:, :-1]
y = df.diagnosis
from sklearn.model_selection import train_test_split
x_train, x_test, y_train, y_test = train_test_split(x, y, test_size= 0.2,
↳random_state= 47)

```

5 Train the Model

```

[ ]: from sklearn.svm import SVC
model = SVC() #Leaving default kernel as kernel="linear" giving both testing_
↳and training accuracy of 100%
clf = model.fit(x_train, y_train)
y_pred = model.predict(x_test)

```

6 Model Evaluation

```
[ ]: from sklearn.metrics import accuracy_score, classification_report
print(f"The training accuracy: {accuracy_score(y_train,model.
      ↪predict(x_train))}\n\nThe testing accuracy: {accuracy_score(y_test,y_pred)}")
print(pd.crosstab(y_test, y_pred))
print(classification_report(y_test, y_pred))
```

The training accuracy: 0.9230769230769231

The testing accuracy: 0.9122807017543859

col_0 0 1

diagnosis

0 67 3

1 7 37

		precision	recall	f1-score	support
	0	0.91	0.96	0.93	70
	1	0.93	0.84	0.88	44
	accuracy			0.91	114
	macro avg	0.92	0.90	0.91	114
	weighted avg	0.91	0.91	0.91	114

6.0.1 There doesn't seem to be any overfitting or underfitting model seems to be performing good with testing accuracy of 91%. With good precision recall and f1-scores

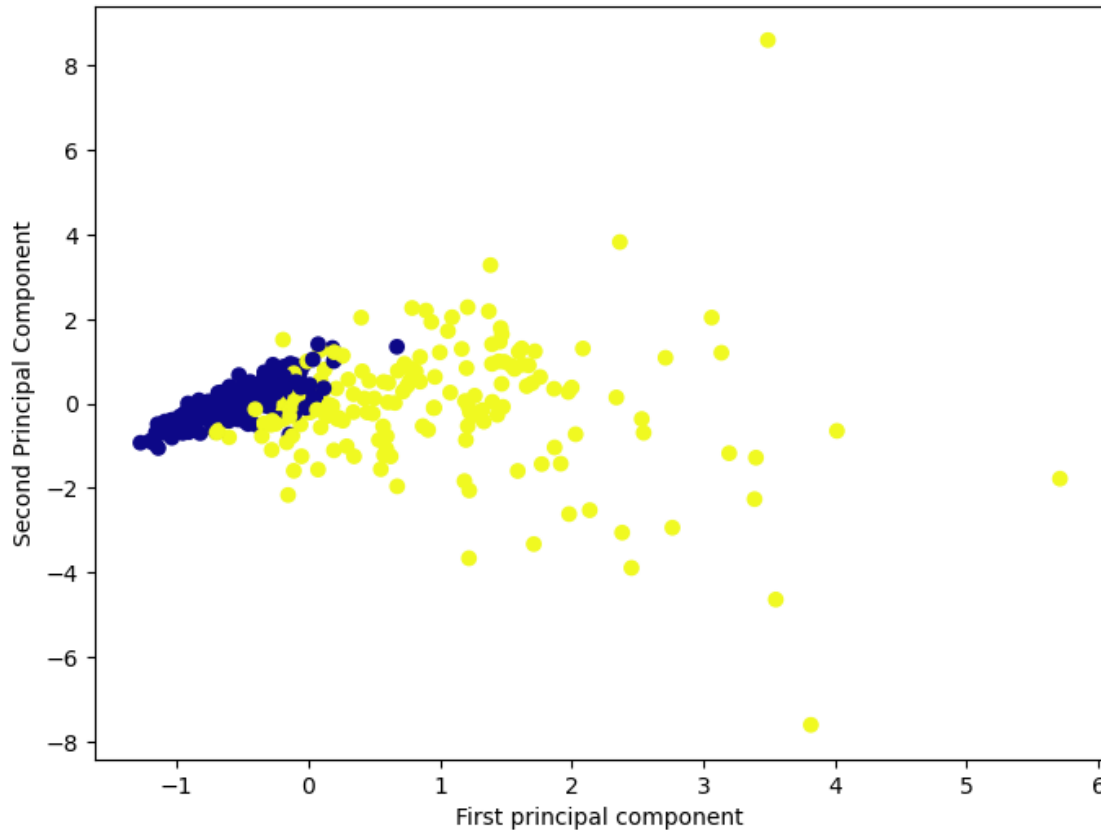
7 Let's Try PCA(Principal Component Analysis)

```
[ ]: from sklearn.decomposition import PCA
from sklearn.preprocessing import StandardScaler
pca = PCA(n_components=4)
scaler = StandardScaler()

X_train = pca.fit_transform(x_train)
X_test = pca.transform(x_test)

X_train = scaler.fit_transform(X_train)
X_test = scaler.transform(X_test)
```

```
[ ]: plt.figure(figsize=(8,6))
plt.scatter(X_train[:,0],X_train[:,1],c=y_train,cmap='plasma')
plt.xlabel('First principal component')
plt.ylabel('Second Principal Component')
plt.show()
```



```
[ ]: from sklearn.svm import SVC
model = SVC() #Leaving default kernel as kernel="linear" giving both testing_
    ↪and training accuracy of 100%
clf = model.fit(X_train, y_train)
y_pred = model.predict(X_test)
```

```
[ ]: from sklearn.metrics import accuracy_score, classification_report
print(f"The training accuracy: {accuracy_score(y_train,model.
    ↪predict(X_train))}\nThe testing accuracy: {accuracy_score(y_test,y_pred)}")
print(pd.crosstab(y_test, y_pred))
print(classification_report(y_test, y_pred))
```

The training accuracy: 0.9648351648351648

The testing accuracy: 0.956140350877193

col_0	0	1
diagnosis		
0	66	4
1	1	43

	precision	recall	f1-score	support
0	0.99	0.94	0.96	70

1	0.91	0.98	0.95	44
accuracy			0.96	114
macro avg	0.95	0.96	0.95	114
weighted avg	0.96	0.96	0.96	114

8 As you we can see PCA increases the model accuracy as there are 34 features to 4 increased the model performace by around 4% which is good indicatoin of dimensionality curse using PCA we solved that

[]:

[]: