

```
In [1]: import pandas as pd
from sklearn.model_selection import train_test_split
from sklearn.preprocessing import StandardScaler
from sklearn.decomposition import PCA
from sklearn.pipeline import Pipeline
from sklearn.ensemble import RandomForestClassifier
from sklearn.metrics import accuracy_score
```

```
In [2]: df=pd.read_csv(r"C:\Users\USER\Downloads\archive\breast-cancer.csv")
df.head()
```

```
Out[2]:
```

	id	diagnosis	radius_mean	texture_mean	perimeter_mean	area_mean	smoothness_mean	compactness_mean	concavity_mean	points_mean	symmetry_mean	fractal_dimension
0	842302	M	17.99	10.38	122.80	1001.0	0.11840	0.27760	0.3001	0.0001	0.0001	0.0001
1	842517	M	20.57	17.77	132.90	1326.0	0.08474	0.07864	0.0869	0.0001	0.0001	0.0001
2	84300903	M	19.69	21.25	130.00	1203.0	0.10960	0.15990	0.1974	0.0001	0.0001	0.0001
3	84348301	M	11.42	20.38	77.58	386.1	0.14250	0.28390	0.2414	0.0001	0.0001	0.0001
4	84358402	M	20.29	14.34	135.10	1297.0	0.10030	0.13280	0.1980	0.0001	0.0001	0.0001

5 rows × 13 columns

```
In [3]: # Preprocess the data
# Encode the diagnosis column (M = 1, B = 0)
df['diagnosis'] = df['diagnosis'].map({'M': 1, 'B': 0})
```

```
In [4]: X=df.drop(columns=['diagnosis','id'])
y=df['diagnosis']
```

```
In [5]: X #Just confirming
```

```
Out[5]:
```

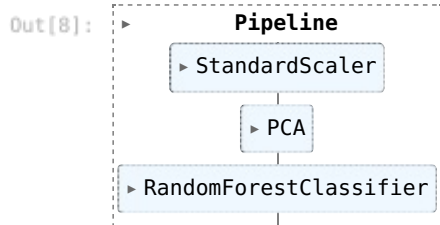
	radius_mean	texture_mean	perimeter_mean	area_mean	smoothness_mean	compactness_mean	concavity_mean	concave points_mean	symmetry_mean	fractal_dimension
0	17.99	10.38	122.80	1001.0	0.11840	0.27760	0.30010	0.14710	0.0001	0.0001
1	20.57	17.77	132.90	1326.0	0.08474	0.07864	0.08690	0.07017	0.0001	0.0001
2	19.69	21.25	130.00	1203.0	0.10960	0.15990	0.19740	0.12790	0.0001	0.0001
3	11.42	20.38	77.58	386.1	0.14250	0.28390	0.24140	0.10520	0.0001	0.0001
4	20.29	14.34	135.10	1297.0	0.10030	0.13280	0.19800	0.10430	0.0001	0.0001
...
564	21.56	22.39	142.00	1479.0	0.11100	0.11590	0.24390	0.13890	0.0001	0.0001
565	20.13	28.25	131.20	1261.0	0.09780	0.10340	0.14400	0.09791	0.0001	0.0001
566	16.60	28.08	108.30	858.1	0.08455	0.10230	0.09251	0.05302	0.0001	0.0001
567	20.60	29.33	140.10	1265.0	0.11780	0.27700	0.35140	0.15200	0.0001	0.0001
568	7.76	24.54	47.92	181.0	0.05263	0.04362	0.00000	0.00000	0.0001	0.0001

569 rows × 11 columns

```
In [6]: X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
```

```
In [7]: # Create a pipeline with a scaler, PCA, and a classifier
pipeline = Pipeline([
    ('scaler', StandardScaler()),          # Step 1: Standardize the data
    ('pca', PCA(n_components=2)),          # Step 2: Apply PCA with 2 components
    ('classifier', RandomForestClassifier(random_state=42)) # Step 3: Fit a classifier (Random Forest)
])
```

```
In [8]: # Fit the pipeline on the training data
pipeline.fit(X_train, y_train)
```



```
In [9]: # Use the pipeline to predict on the test data
y_pred = pipeline.predict(X_test)
```

```
In [10]: # Evaluate the model
accuracy = accuracy_score(y_test, y_pred)
print(f'Accuracy: {accuracy}')
```

Accuracy: 0.9824561403508771

Loading [MathJax]/jax/output/CommonHTML/fonts/TeX/fontdata.js