# Forward School

**Program Code: J620-002-4:2020**

**Program Name: FRONT-END SOFTWARE DEVELOPMENT**

**Title : Exe22 - Bagging and Boosting Exercise**

**Name: Chuay Xiang Ze**

**IC Number: 021224070255**

**Date : 20/07/2023**

**Introduction : Learning how to use bagging and boosting.**

**Conclusion : Managed to complete tasks relating to the topic.**

## Bagging and Boosting Exercise

Reference: ([https://www.datacamp.com/community/tutorials/ensemble-learning-python (https://www.datacamp.com/community/tutorials/ensemble-learning-python)](https://www.datacamp.com/community/tutorials/ensemble-learning-python))

## Bagging Method

In [20]:

```python
import pandas as pd
import numpy as np
from sklearn.impute import SimpleImputer
imputer = SimpleImputer(missing_values=np.nan, strategy='mean')
from sklearn.preprocessing import MinMaxScaler
```

```python
import pandas as pd
data = pd.read_csv('https://archive.ics.uci.edu/ml/machine-learning-databases/breast-can
                   header=None)
data.columns = ['Sample code', 'Clump Thickness', 'Uniformity of Cell Size', 'Uniformity
                'Marginal Adhesion', 'Single Epithelial Cell Size', 'Bare Nuclei', 'Blan
                'Normal Nucleoli', 'Mitoses','Class']

data = data.drop(['Sample code'],axis=1)
print('Number of instances = %d' % (data.shape[0]))
print('Number of attributes = %d' % (data.shape[1]))
data.head()
```

```
Number of instances = 699
Number of attributes = 10
```

Out[21]:

| | Clump Thickness | Uniformity of Cell Size | Uniformity of Cell Shape | Marginal Adhesion | Single Epithelial Cell Size | Bare Nuclei | Bland Chromatin | Normal Nucleoli | Mitos |
|---|---|---|---|---|---|---|---|---|---|
| 0 | 5 | 1 | 1 | 1 | 2 | 1 | 3 | 1 | |
| 1 | 5 | 4 | 4 | 5 | 7 | 10 | 3 | 2 | |
| 2 | 3 | 1 | 1 | 1 | 2 | 2 | 3 | 1 | |
| 3 | 6 | 8 | 8 | 1 | 3 | 4 | 3 | 7 | |
| 4 | 4 | 1 | 1 | 3 | 2 | 1 | 3 | 1 | |

In [22]:

```python
data.describe()
```

Out[22]:

| | Clump Thickness | Uniformity of Cell Size | Uniformity of Cell Shape | Marginal Adhesion | Single Epithelial Cell Size | Bland Chromatin | Normal Nucleoli |
|---|---|---|---|---|---|---|---|
| count | 699.000000 | 699.000000 | 699.000000 | 699.000000 | 699.000000 | 699.000000 | 699.000000 |
| mean | 4.417740 | 3.134478 | 3.207439 | 2.806867 | 3.216023 | 3.437768 | 2.866953 |
| std | 2.815741 | 3.051459 | 2.971913 | 2.855379 | 2.214300 | 2.438364 | 3.053634 |
| min | 1.000000 | 1.000000 | 1.000000 | 1.000000 | 1.000000 | 1.000000 | 1.000000 |
| 25% | 2.000000 | 1.000000 | 1.000000 | 1.000000 | 2.000000 | 2.000000 | 1.000000 |
| 50% | 4.000000 | 1.000000 | 1.000000 | 1.000000 | 2.000000 | 3.000000 | 1.000000 |
| 75% | 6.000000 | 5.000000 | 5.000000 | 4.000000 | 4.000000 | 5.000000 | 4.000000 |
| max | 10.000000 | 10.000000 | 10.000000 | 10.000000 | 10.000000 | 10.000000 | 10.000000 |

In [23]:

```python
data.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 699 entries, 0 to 698
Data columns (total 10 columns):
 #   Column                       Non-Null Count  Dtype
---  ------                       --------------  -----
 0   Clump Thickness              699 non-null    int64
 1   Uniformity of Cell Size      699 non-null    int64
 2   Uniformity of Cell Shape     699 non-null    int64
 3   Marginal Adhesion            699 non-null    int64
 4   Single Epithelial Cell Size  699 non-null    int64
 5   Bare Nuclei                  699 non-null    object
 6   Bland Chromatin              699 non-null    int64
 7   Normal Nucleoli              699 non-null    int64
 8   Mitoses                      699 non-null    int64
 9   Class                        699 non-null    int64
dtypes: int64(9), object(1)
memory usage: 54.7+ KB
```

In [24]:

```python
data['Bare Nuclei']
```

Out[24]:

```
0       1
1      10
2       2
3       4
4       1
       ..
694     2
695     1
696     3
697     4
698     5
Name: Bare Nuclei, Length: 699, dtype: object
```

In [25]:

```python
data.replace('?',0, inplace=True)
data['Bare Nuclei']
```

Out[25]:

```
0       1
1      10
2       2
3       4
4       1
       ..
694     2
695     1
696     3
697     4
698     5
Name: Bare Nuclei, Length: 699, dtype: object
```

In [26]:

```python
# Convert the DataFrame object into NumPy array otherwise you will not be able to impute
values = data.values

# Now impute it

imputedData = imputer.fit_transform(values)
```

In [27]:

```python
scaler = MinMaxScaler(feature_range=(0, 1))
normalizedData = scaler.fit_transform(imputedData)
```

In [28]:

```python
# Bagged Decision Trees for Classification - necessary dependencies

# test classification dataset
from sklearn.datasets import make_classification
# define dataset
X, y = make_classification(n_samples=1000, n_features=20, n_informative=15, n_redundant=
# summarize the dataset
print(X.shape, y.shape)
# evaluate bagging algorithm for classification
from numpy import mean
from numpy import std
from sklearn.model_selection import cross_val_score
from sklearn.model_selection import RepeatedStratifiedKFold
from sklearn.ensemble import BaggingClassifier
# define dataset
X, y = make_classification(n_samples=1000, n_features=20, n_informative=15, n_redundant=
# define the model
model = BaggingClassifier()
# evaluate the model
cv = RepeatedStratifiedKFold(n_splits=10, n_repeats=3, random_state=1)
n_scores = cross_val_score(model, X, y, scoring='accuracy', cv=cv, n_jobs=-1, error_scor
# report performance
print('Accuracy: %.3f (%.3f)' % (mean(n_scores), std(n_scores)))
```

```
(1000, 20) (1000,)
Accuracy: 0.866 (0.031)
```

In [31]:

```python
# Segregate the features from the labels
X = data.drop('Class', axis = 1)
y = data['Class']
```

In [41]:

```python
from sklearn.ensemble import BaggingRegressor
from sklearn.datasets import make_regression
from sklearn.model_selection import RepeatedKFold

# evaluate the model
# evaluate the model
cv = RepeatedStratifiedKFold(n_splits=10, n_repeats=3, random_state=1)
n_scores = cross_val_score(model, X, y, scoring='accuracy', cv=cv, n_jobs=-1, error_scor
print('Accuracy: %.3f (%.3f)' % (mean(n_scores), std(n_scores)))
```

Accuracy: 0.957 (0.021)

In [11]:

```
D:\Anaconda3\lib\site-packages\sklearn\model_selection\_split.py:296: Futu
reWarning: Setting a random_state has no effect since shuffle is False. Th
is will raise an error in 0.24. You should leave random_state to its defau
lt (None), or set shuffle=True.
  FutureWarning
```

0.9585714285714285

# Boosting Method

In [49]:

```python
from sklearn.ensemble import AdaBoostClassifier
from sklearn import model_selection
seed = 7
num_trees = 70
kfold = model_selection.KFold(n_splits=10, random_state=seed, shuffle=True)
model = AdaBoostClassifier(n_estimators=num_trees, random_state=seed)
results = model_selection.cross_val_score(model, X, y, cv=kfold)
print(results.mean())
```

0.9599378881987578

# Exercise 1 Perform classification using the Titanic dataset using the classifiers that you already know (Dtree and RF)

```python
#Preprocessing the entire Titanic dataset
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
%matplotlib inline

t_dataset = pd.read_csv('./titanic.csv')
t_dataset
```

| | Survived | Pclass | Name | Sex | Age | Siblings/Spouses Aboard | Parents/Children Aboard | Fare |
|---|---|---|---|---|---|---|---|---|
| **0** | 0 | 3 | Mr. Owen Harris Braund | male | 22.0 | 1 | 0 | 7.2500 |
| **1** | 1 | 1 | Mrs. John Bradley (Florence Briggs Thayer) Cum... | female | 38.0 | 1 | 0 | 71.2833 |
| **2** | 1 | 3 | Miss. Laina Heikkinen | female | 26.0 | 0 | 0 | 7.9250 |
| **3** | 1 | 1 | Mrs. Jacques Heath (Lily May Peel) Futrelle | female | 35.0 | 1 | 0 | 53.1000 |
| **4** | 0 | 3 | Mr. William Henry Allen | male | 35.0 | 0 | 0 | 8.0500 |
| **...** | ... | ... | ... | ... | ... | ... | ... | ... |
| **882** | 0 | 2 | Rev. Juozas Montvila | male | 27.0 | 0 | 0 | 13.0000 |
| **883** | 1 | 1 | Miss. Margaret Edith Graham | female | 19.0 | 0 | 0 | 30.0000 |
| **884** | 0 | 3 | Miss. Catherine Helen Johnston | female | 7.0 | 1 | 2 | 23.4500 |
| **885** | 1 | 1 | Mr. Karl Howell Behr | male | 26.0 | 0 | 0 | 30.0000 |
| **886** | 0 | 3 | Mr. Patrick Dooley | male | 32.0 | 0 | 0 | 7.7500 |

887 rows × 8 columns

```
#drop name column
t_dataset = t_dataset.drop('Name', axis = 1)
t_dataset
```

Out[51]:

| | Survived | Pclass | Sex | Age | Siblings/Spouses Aboard | Parents/Children Aboard | Fare |
|---|---|---|---|---|---|---|---|
| **0** | 0 | 3 | male | 22.0 | 1 | 0 | 7.2500 |
| **1** | 1 | 1 | female | 38.0 | 1 | 0 | 71.2833 |
| **2** | 1 | 3 | female | 26.0 | 0 | 0 | 7.9250 |
| **3** | 1 | 1 | female | 35.0 | 1 | 0 | 53.1000 |
| **4** | 0 | 3 | male | 35.0 | 0 | 0 | 8.0500 |
| **...** | ... | ... | ... | ... | ... | ... | ... |
| **882** | 0 | 2 | male | 27.0 | 0 | 0 | 13.0000 |
| **883** | 1 | 1 | female | 19.0 | 0 | 0 | 30.0000 |
| **884** | 0 | 3 | female | 7.0 | 1 | 2 | 23.4500 |
| **885** | 1 | 1 | male | 26.0 | 0 | 0 | 30.0000 |
| **886** | 0 | 3 | male | 32.0 | 0 | 0 | 7.7500 |

887 rows × 7 columns

```
#encode categorical data into numerical value
from sklearn import preprocessing
reset = {"male": 1, "female": 0}
t_dataset = t_dataset.replace({"Sex": reset})
t_dataset = t_dataset.dropna()
t_dataset
```

| | Survived | Pclass | Sex | Age | Siblings/Spouses Aboard | Parents/Children Aboard | Fare |
|---|---|---|---|---|---|---|---|
| 0 | 0 | 3 | 1 | 22.0 | 1 | 0 | 7.2500 |
| 1 | 1 | 1 | 0 | 38.0 | 1 | 0 | 71.2833 |
| 2 | 1 | 3 | 0 | 26.0 | 0 | 0 | 7.9250 |
| 3 | 1 | 1 | 0 | 35.0 | 1 | 0 | 53.1000 |
| 4 | 0 | 3 | 1 | 35.0 | 0 | 0 | 8.0500 |
| ... | ... | ... | ... | ... | ... | ... | ... |
| 882 | 0 | 2 | 1 | 27.0 | 0 | 0 | 13.0000 |
| 883 | 1 | 1 | 0 | 19.0 | 0 | 0 | 30.0000 |
| 884 | 0 | 3 | 0 | 7.0 | 1 | 2 | 23.4500 |
| 885 | 1 | 1 | 1 | 26.0 | 0 | 0 | 30.0000 |
| 886 | 0 | 3 | 1 | 32.0 | 0 | 0 | 7.7500 |

887 rows × 7 columns

```
#create a copy of the cleaned dataset
t_datacopy = t_dataset.copy()
t_datacopy
```

| | Survived | Pclass | Sex | Age | Siblings/Spouses Aboard | Parents/Children Aboard | Fare |
|---|---|---|---|---|---|---|---|
| **0** | 0 | 3 | 1 | 22.0 | 1 | 0 | 7.2500 |
| **1** | 1 | 1 | 0 | 38.0 | 1 | 0 | 71.2833 |
| **2** | 1 | 3 | 0 | 26.0 | 0 | 0 | 7.9250 |
| **3** | 1 | 1 | 0 | 35.0 | 1 | 0 | 53.1000 |
| **4** | 0 | 3 | 1 | 35.0 | 0 | 0 | 8.0500 |
| **...** | ... | ... | ... | ... | ... | ... | ... |
| **882** | 0 | 2 | 1 | 27.0 | 0 | 0 | 13.0000 |
| **883** | 1 | 1 | 0 | 19.0 | 0 | 0 | 30.0000 |
| **884** | 0 | 3 | 0 | 7.0 | 1 | 2 | 23.4500 |
| **885** | 1 | 1 | 1 | 26.0 | 0 | 0 | 30.0000 |
| **886** | 0 | 3 | 1 | 32.0 | 0 | 0 | 7.7500 |

887 rows × 7 columns

In [93]:

```python
#define dependent variable and independent variable
X = t_datacopy.drop('Survived', axis = 1)
y = t_datacopy['Survived']
X_list = X.values.tolist()[0]
y_list = y.tolist()

print(X_list)
print(y_list)
```

```
[3.0, 1.0, 22.0, 1.0, 0.0, 7.25]
[0, 1, 1, 1, 0, 0, 0, 0, 1, 1, 1, 1, 0, 0, 0, 1, 0, 1, 0, 1, 0, 1, 1, 1,
 0, 1, 0, 0, 1, 0, 0, 1, 1, 0, 0, 0, 1, 0, 0, 1, 0, 0, 1, 1, 0, 0, 1, 0, 0,
 0, 0, 1, 1, 0, 1, 1, 0, 1, 0, 0, 1, 0, 0, 0, 1, 1, 0, 1, 0, 0, 0, 0, 0, 1,
 0, 0, 0, 1, 1, 0, 1, 1, 0, 1, 1, 0, 0, 1, 0, 0, 0, 0, 0, 0, 0, 0, 1, 1, 0,
 0, 0, 0, 0, 0, 0, 1, 1, 0, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 0,
 1, 0, 1, 1, 0, 0, 0, 0, 1, 0, 0, 1, 0, 0, 0, 0, 1, 1, 0, 0, 0, 1, 0, 0, 0,
 0, 1, 0, 0, 0, 0, 1, 0, 0, 0, 0, 1, 0, 0, 0, 1, 1, 0, 0, 0, 0, 0, 1, 0, 0,
 0, 0, 0, 0, 0, 0, 0, 0, 1, 1, 0, 1, 1, 0, 0, 1, 0, 1, 1, 1, 1, 0, 0, 1, 0,
 0, 0, 0, 0, 1, 0, 0, 1, 1, 1, 0, 1, 0, 0, 0, 1, 1, 0, 1, 0, 1, 0, 0, 0, 1,
 0, 1, 0, 0, 0, 1, 0, 0, 1, 0, 0, 0, 1, 0, 0, 0, 1, 0, 0, 0, 0, 0, 1, 1, 0,
 0, 0, 0, 0, 0, 1, 1, 1, 1, 1, 1, 0, 0, 0, 0, 0, 1, 1, 1, 0, 1, 1, 0, 1, 1,
 0, 0, 0, 1, 0, 0, 0, 1, 0, 0, 1, 0, 1, 1, 1, 1, 0, 0, 0, 0, 0, 0, 1, 1, 1,
 1, 0, 1, 0, 1, 1, 1, 0, 1, 1, 1, 0, 0, 0, 1, 1, 0, 1, 1, 0, 0, 1, 1, 0, 1,
 0, 1, 1, 1, 1, 0, 0, 0, 1, 0, 0, 1, 1, 0, 1, 1, 0, 0, 0, 1, 1, 1, 1, 0, 0,
 0, 0, 0, 0, 0, 1, 0, 1, 1, 0, 0, 0, 0, 0, 0, 1, 1, 1, 1, 1, 0, 0, 0, 0, 1,
 1, 0, 0, 0, 1, 1, 0, 1, 0, 0, 0, 1, 0, 1, 1, 1, 0, 1, 1, 0, 0, 0, 0, 1, 1,
 0, 0, 0, 0, 0, 0, 1, 0, 0, 0, 0, 1, 0, 1, 0, 1, 1, 0, 0, 0, 0, 0, 0, 0, 1,
 1, 0, 1, 1, 1, 1, 0, 0, 1, 0, 1, 0, 0, 1, 0, 0, 1, 1, 1, 1, 1, 1, 1, 0, 0,
 0, 1, 0, 1, 0, 1, 1, 0, 1, 0, 0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 1, 1, 0, 0, 0,
 0, 0, 1, 0, 0, 0, 1, 1, 0, 1, 0, 0, 1, 0, 0, 0, 0, 0, 0, 1, 0, 0, 0, 0, 0,
 0, 0, 1, 0, 1, 1, 0, 1, 1, 0, 1, 1, 0, 0, 1, 0, 1, 0, 1, 0, 0, 1, 0, 0, 1,
 0, 0, 0, 1, 0, 0, 1, 0, 1, 0, 1, 0, 1, 1, 0, 0, 1, 0, 0, 1, 1, 0, 1, 1, 0,
 0, 1, 1, 0, 1, 0, 1, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 1, 1, 1, 1, 0, 0, 1,
 1, 0, 1, 1, 1, 0, 0, 0, 1, 0, 1, 0, 0, 0, 1, 0, 0, 0, 0, 1, 0, 0, 1, 1, 0,
 0, 0, 1, 0, 0, 1, 1, 1, 0, 0, 1, 0, 0, 1, 0, 0, 1, 0, 0, 1, 1, 0, 0, 0, 0,
 1, 0, 0, 1, 0, 1, 0, 0, 1, 0, 0, 0, 0, 0, 1, 0, 1, 1, 1, 0, 1, 0, 1, 0, 1,
 0, 1, 0, 0, 0, 0, 0, 0, 1, 0, 0, 0, 1, 0, 0, 0, 0, 1, 1, 0, 0, 1, 0, 0, 0,
 1, 0, 1, 0, 1, 0, 0, 0, 0, 0, 0, 0, 1, 1, 1, 1, 0, 0, 0, 0, 1, 0, 0, 1, 1,
 0, 0, 0, 0, 1, 1, 1, 1, 1, 0, 1, 0, 0, 0, 1, 1, 0, 1, 0, 0, 0, 1, 0, 1, 1,
 0, 0, 1, 0, 0, 0, 0, 0, 0, 1, 0, 0, 1, 0, 1, 0, 1, 0, 0, 1, 0, 0, 1, 1, 0,
 0, 1, 1, 0, 0, 0, 1, 0, 0, 1, 1, 0, 1, 0, 0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 1,
 0, 1, 1, 1, 0, 0, 0, 0, 1, 0, 1, 0, 0, 0, 0, 0, 0, 0, 1, 1, 0, 0, 0, 1, 1,
 1, 1, 0, 0, 0, 0, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 1, 0, 1, 0, 0, 0, 1,
 1, 1, 1, 1, 0, 0, 0, 1, 0, 0, 1, 1, 0, 0, 1, 0, 0, 0, 0, 0, 0, 1, 0, 0, 0,
 1, 0, 1, 1, 1, 1, 0, 0, 0, 1, 0, 0, 1, 1, 0, 0, 1, 0, 1, 0, 0, 1, 1, 0, 0,
 0, 1, 1, 0, 0, 0, 0, 0, 0, 1, 0, 1, 0]
```

In [86]:

```python
#Split the dataset into the Training and the Test set. Set the test set to 0.3
from sklearn.model_selection import train_test_split
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3, random_state=42
```

In [87]:

```python
from sklearn.tree import DecisionTreeClassifier
from sklearn.metrics import accuracy_score
```

In [107]:

```python
#Decision Tree object
clf = DecisionTreeClassifier()

# Train Decision Tree Classifer
clf.fit(X_train,y_train)

#Predict the response for test dataset
y_pred_dt = clf.predict(X_test)

# Model Accuracy, how often is the classifier correct
accuracy_dt = accuracy_score(y_test, y_pred_dt)
print(y_pred_dt)
print("Accuracy:", accuracy_dt)
```

```
[1 0 0 0 0 1 0 0 0 0 1 1 0 0 0 0 0 1 0 0 0 1 0 0 0 0 0 0 1 0 0 0 0 1 1 0 0
 1 1 1 1 0 1 0 0 0 1 0 0 0 0 0 0 0 1 0 1 1 1 1 0 1 0 1 1 0 0 1 0 0 1 1 1 0
 1 1 1 1 1 0 0 0 0 0 1 1 0 0 1 0 1 1 0 1 0 0 0 0 1 0 0 0 1 1 0 0 0 1 0 0 1
 0 0 0 0 0 1 0 0 0 1 0 1 0 0 0 0 1 1 1 0 1 0 0 0 0 0 0 0 1 0 0 0 0 1 0 0
 0 0 0 1 1 0 0 1 0 1 1 0 0 1 0 1 0 1 0 1 0 1 1 1 1 0 1 1 1 0 0 0 0 1 1 0 1
 1 0 0 0 0 0 1 0 0 1 0 0 0 1 0 1 1 1 0 0 0 0 0 1 0 1 0 0 1 1 0 1 0 1 0 0 0
 0 1 0 0 0 0 0 0 0 1 0 1 1 1 1 1 0 0 1 0 1 1 0 1 0 0 1 1 0 0 1 0 0 0 1 0 0
 1 0 0 1 0 0 1 0]
Accuracy: 0.7790262172284644
```

In [ ]:

```python
#Random forest
from sklearn.ensemble import RandomForestClassifier

# Create the model with 100 trees
num_trees = 100
model = RandomForestClassifier(n_estimators=num_trees)

# Fit on training data
model.fit(X, y)

# Probabilities for each class
y_pred_rf = model.predict(X_test)
y_prob_rf = model.predict_proba(X_test)
# Assuming y_prob_rf contains the probabilities from model.predict_proba(X_test)

# Flatten the array
y_prob_rf_flattened = y_prob_rf.flatten()

print(y_pred_rf)
# Print the flattened array
print(y_prob_rf_flattened)

accuracy_rf = accuracy_score(y_test, y_pred_rf)
print("Accuracy:", accuracy_rf)
```

```
[1 0 0 1 0 1 0 0 1 1 1 1 1 0 1 0 0 1 0 0 0 1 0 1 0 0 0 0 0 1 0 0 1 1 1 1 0 0
 1 1 0 1 0 1 0 0 1 0 0 0 0 0 1 0 0 1 0 1 1 0 1 0 1 0 1 1 0 0 1 0 0 0 1 0 0
 0 1 1 0 0 0 0 0 0 1 1 0 0 0 0 1 0 0 1 0 0 0 0 1 1 0 0 0 1 0 0 0 0 0 1 1
 0 0 0 0 0 1 1 0 0 0 1 0 0 0 0 1 0 1 1 0 1 0 0 0 0 0 0 0 0 0 0 0 0 1 0 0
 0 0 0 1 0 0 0 1 0 1 1 0 0 0 1 1 0 1 0 1 0 1 0 1 1 1 0 1 1 1 1 0 0 0 0 1 1 0 1
 0 0 0 0 1 0 1 0 0 0 0 0 0 1 1 0 1 1 0 1 0 0 0 0 1 0 0 0 0 1 1 0 0 0 1 0 0 0
 0 0 0 1 0 0 1 1 0 0 1 1 0 0 1 1 0 0 1 0 1 0 0 1 0 0 1 0 0 1 1 0 0 1 0 0 0 1 0 0
 1 0 0 1 1 1 1 0]
[0.27       0.73       0.93       0.07       0.99       0.01
 0.29       0.71       0.98       0.02       0.3        0.7
 1.         0.         1.         0.         0.09       0.91
 0.26       0.74       0.045      0.955      0.03       0.97
 0.12       0.88       0.97       0.03       0.26       0.74
 0.96083333 0.03916667 0.99       0.01       0.         1.
 0.84875    0.15125    0.9855     0.0145     0.98       0.02
 0.01       0.99       0.87       0.13       0.26       0.74
 0.91       0.09       0.848      0.152      0.99       0.01
 0.98       0.02       0.97       0.03       0.34       0.66
 0.74       0.26       0.9        0.1        0.17       0.83
 0.21       0.79       0.1        0.9        0.94       0.06
 0.96       0.04       0.025      0.975      0.49942857 0.50057143
 0.72       0.28       0.03       0.97       1.         0.
 0.00333333 0.99666667 0.98666667 0.01333333 0.98       0.02
 0.319      0.681      0.99       0.01       1.         0.
 0.98       0.02       1.         0.         0.98       0.02
 0.32       0.68       0.9        0.1        1.         0.
 0.02       0.98       0.97       0.03       0.2425     0.7575
 0.01       0.99       0.76       0.24       0.26       0.74
 0.98       0.02       0.06       0.94       0.66897619 0.33102381
 0.01       0.99       0.37       0.63       0.97       0.03
 1.         0.         0.04       0.96       0.99       0.01
 1.         0.         0.62142063 0.37857937 0.         1.
 0.50238095 0.49761905 1.         0.         0.82364286 0.17635714
 0.         1.         0.03       0.97       0.76       0.24
 0.67       0.33       1.         0.         1.         0.
 0.91666667 0.08333333 0.99       0.01       0.98       0.02
 0.02       0.98       0.         1.         0.78       0.22
 0.98       0.02       0.62142063 0.37857937 0.85483333 0.14516667
 0.12       0.88       0.72       0.28       0.985      0.015
 0.03       0.97       0.97       0.03       0.8        0.2
 0.99666667 0.00333333 0.99       0.01       0.03       0.97
 0.06       0.94       0.87       0.13       0.55816667 0.44183333
 0.87333333 0.12666667 0.01       0.99       1.         0.
 1.         0.         0.92       0.08       0.76       0.24
 1.         0.         0.24       0.76       0.03333333 0.96666667
 1.         0.         0.97166667 0.02833333 0.78       0.22
 0.95166667 0.04833333 0.55516667 0.44483333 0.22       0.78
 0.16       0.84       0.94       0.06       0.99       0.01
 0.72       0.28       0.2525     0.7475     0.83       0.17
 1.         0.         0.82       0.18       0.97       0.03
 0.37       0.63       0.63       0.37       0.03       0.97
 0.03       0.97       0.92       0.08       0.18       0.82
 1.         0.         1.         0.         1.         0.
 0.87       0.13       0.99       0.01       1.         0.
 0.96       0.04       1.         0.         0.69       0.31
 0.98       0.02       0.68       0.32       0.87       0.13
 0.73       0.27       0.01       0.99       0.99       0.01
 0.92683333 0.07316667 0.81       0.19       1.         0.
 0.98       0.02       0.09       0.91       0.965      0.035
 0.97       0.03       0.86       0.14       0.08166667 0.91833333
 1.         0.         0.04       0.96       0.06       0.94
```

```
 0.99        0.01        1.          0.          0.76        0.24
 0.3         0.7         0.23        0.77        0.99        0.01
 0.12        0.88        0.95        0.05        0.03        0.97
 0.97        0.03        0.01        0.99        0.01        0.99
 0.13        0.87        0.83233333 0.16766667 0.41         0.59
 0.19        0.81        0.07        0.93        0.07        0.93
 0.88        0.12        0.74323016 0.25676984 0.94333333 0.05666667
 0.87        0.13        0.01333333 0.98666667 0.02         0.98
 1.          0.          0.07333333 0.92666667 0.77         0.23
 1.          0.          1.          0.          0.91        0.09
 0.28        0.72        1.          0.          0.01        0.99
 1.          0.          0.55516667 0.44483333 0.66         0.34
 1.          0.          0.69        0.31        0.41        0.59
 0.          1.          0.95        0.05        0.08        0.92
 0.          1.          0.76        0.24        0.32        0.68
 0.98        0.02        0.94        0.06        0.99        0.01
 1.          0.          0.          1.          0.96        0.04
 0.85        0.15        0.9         0.1         0.96        0.04
 0.018       0.982       0.24        0.76        1.          0.
 0.82        0.18        0.92        0.08        0.01        0.99
 0.98        0.02        0.95        0.05        1.          0.
 1.          0.          0.62142063 0.37857937 1.          0.
 0.22        0.78        0.99        0.01        1.          0.
 0.49128571 0.50871429 0.3         0.7         0.93        0.07
 0.89666667 0.10333333 0.23        0.77        0.35        0.65
 0.70933333 0.29066667 0.88        0.12        0.11128571 0.88871429
 0.12        0.88        0.99        0.01        0.95        0.05
 0.15        0.85        0.965       0.035       0.17        0.83
 0.77        0.23        0.98        0.02        0.06        0.94
 1.          0.          0.55816667 0.44183333 0.03         0.97
 0.          1.          0.9625      0.0375      0.98416667 0.01583333
 0.02        0.98        0.99        0.01        1.          0.
 1.          0.          0.          1.          0.72        0.28
 0.93        0.07        0.18        0.82        0.99        0.01
 0.75        0.25        0.2         0.8         0.38697619 0.61302381
 0.15        0.85        0.          1.          1.          0.         ]
Accuracy: 0.9812734082397003
```

# Exercise 2 Perform classification using the Titanic dataset using the classifiers that you already know and with feature selection and dimension reduction. Which gives you the best result?

In [ ]:

In [122]:

```python
# Step 1: Import the required libraries
import pandas as pd
from sklearn.model_selection import train_test_split
from sklearn.preprocessing import StandardScaler
from sklearn.decomposition import PCA
from sklearn.tree import DecisionTreeClassifier
from sklearn.metrics import accuracy_score


# Step 2: Feature Scaling with StandardScaler
sc = StandardScaler()
X_train_scaled = sc.fit_transform(X_train)
X_test_scaled = sc.transform(X_test)

# Step 3: Apply PCA
pca = PCA(n_components=2)
X_train_pca = pca.fit_transform(X_train_scaled)
X_test_pca = pca.transform(X_test_scaled)

# Step 4: Create and train Decision Tree Classifier
clf = DecisionTreeClassifier()
clf.fit(X_train_pca, y_train)

# Step 5: Predict the response for test dataset
y_pred = clf.predict(X_test_pca)
explained_variance_ratio = pca.explained_variance_ratio_

print(explained_variance_ratio)
print(y_pred)
accuracy = accuracy_score(y_test, y_pred)
print(accuracy)
```

```
[0.30149571 0.29208691]
[0 0 0 0 0 1 1 0 0 1 1 1 1 0 0 0 0 0 1 0 0 1 1 0 0 0 0 0 0 1 1 1 0 1 1 0 0
 1 1 1 1 0 1 0 0 0 0 0 0 0 0 0 0 0 0 0 1 1 0 0 0 0 1 0 0 0 0 0 0 0 1 1 0
 0 0 0 1 1 0 0 0 0 0 1 1 1 0 0 0 1 1 0 0 0 1 0 0 1 1 1 0 0 1 0 0 1 1 0 0 1
 0 0 0 0 0 0 1 0 0 1 0 0 0 1 1 0 1 1 1 1 0 0 0 1 0 0 0 0 0 1 0 0 0 1 0 1
 1 0 0 0 1 0 0 1 0 1 0 1 0 1 0 1 0 1 0 1 1 0 0 0 0 1 1 0 0 0 1 1 1 0 1
 1 0 0 1 0 0 1 0 0 1 0 0 1 1 1 1 1 1 0 0 0 0 0 1 0 1 0 0 1 1 0 1 0 1 1 0 0
 0 0 1 1 0 0 0 0 0 0 1 1 1 0 0 1 0 0 1 0 1 1 0 1 0 0 0 1 0 0 1 0 0 0 1 0 0
 0 0 0 1 0 1 0 0]
0.6928838951310862
```

```python
#StandardScaler
from sklearn.preprocessing import StandardScaler
sc = StandardScaler()


#PCA & Pick up from the point where dataset has been split
from sklearn.decomposition import PCA


#Decision Tree object


# Train Decision Tree Classifer


#Predict the response for test dataset


# Model Accuracy, how often is the classifier correct?
```

```
[0.30978712 0.28517963]
[1 0 0 1 0 0 1 1 0 1 1 0 0 0 0 1 0 0 1 0 0 1 1 1 1 0 1 0 0 0 1 1 1 0 1 0
 0 0 0 0 1 0 0 1 1 0 0 1 0 0 1 0 0 0 1 0 0 0 1 0 0 0 0 0 0 0 0 0 1 0 0 1
 0 1 0 0 0 0 0 1 1 1 0 1 0 1 1 0 1 0 0 0 0 1 0 0 0 0 1 1 1 0 0 1 1 0 1 1
 0 1 0 1 0 0 1 0 0 1 0 0 1 1 0 1 0 1 0 1 1 0 1 0 0 0 0 0 1 0 1 0 0 1 0 1 1
 0 0 0 0 1 1 1 1 1 1 0 0 1 1 1 0 1 1 0 0 0 1 0 1 0 1 0 0 0 1 0 1 1 0 0 0 0
 0 0 0 1 1 0 0 0 1 0 0 1 0 0 1 0 0 0 1 1 0 0 1 0 0 1 1 0 1 1 0 0 1 0 1 0 1
 1 0 0 1 0 1 0 1 0 1 0 1 0 1 0 0 1 0 0 1 1 1 0 0 0 1 0 1 0 1 1 0 0 0 1 0 1 0 0
 1 1 1 0 1 1 0 1]
Accuracy: 0.7191011235955056
```

In [114]:

```python
#rebuild analytical dataset & create a copy of the cleaned dataset

#define dependent variable and independent variable


#Split the dataset into the Training and the Test set. Set the test set to 0.3
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3, random_state=42
```

```python
from sklearn.feature_selection import SelectFromModel
num_trees = 10000
model = RandomForestClassifier(n_estimators=num_trees)

# Create a selector object that will use the random forest classifier to identify
# features that have an importance of more than 0.15
feature_selector = SelectFromModel(estimator=model, threshold=0.15)

# Train the selector
feature_selector.fit(X_train, y_train)

# Transform the data to create a new dataset containing only the most important features
X_train_selected = feature_selector.transform(X_train)
X_test_selected = feature_selector.transform(X_test)

# Create a new random forest classifier for the most important features
model_selected = RandomForestClassifier(n_estimators=num_trees)

# Train the new classifier on the new dataset containing the most important features
model_selected.fit(X_train_selected, y_train)

# Apply the limited classifier to the test data
y_pred_selected = model_selected.predict(X_test_selected)

# View the accuracy of our limited feature (selected features) model
accuracy_selected = accuracy_score(y_test, y_pred_selected)
print("Accuracy with selected features:", accuracy_selected)
```

Accuracy with selected features: 0.7715355805243446

```python
#RF Feature Selector
import numpy as np
from sklearn.ensemble import RandomForestClassifier
from sklearn import datasets
from sklearn.model_selection import train_test_split
from sklearn.feature_selection import SelectFromModel
from sklearn.metrics import accuracy_score

# Create a random forest classifier (10000 trees)
num_trees = 10000

# Train the classifier
model = RandomForestClassifier()

# Create a selector object that will use the random forest classifier to identify
# features that have an importance of more than 0.15


# Train the selector


# Transform the data to create a new dataset containing only the most important features
# Note: We have to apply the transform to both the training X and test X data.


# Create a new random forest classifier for the most important features


# Train the new classifier on the new dataset containing the most important features


# Apply The Limited Classifier To The Test Data


# View The Accuracy Of Our Limited Feature (2 Features) Model
```

0.7790262172284644

# Exercise 3 Perform classification using the Titanic dataset using bagging and boosting (choose 1 bagging and 1 boosting algo)

```
!pip install xgboost
```

```
Collecting xgboost
  Downloading xgboost-1.7.6-py3-none-win_amd64.whl (70.9 MB)
                                    0.0/70.9 MB ? eta -:--:--
                                    0.2/70.9 MB 4.8 MB/s eta
0:00:15
                                    0.4/70.9 MB 6.0 MB/s eta
0:00:12
                                    0.7/70.9 MB 4.3 MB/s eta
0:00:17
                                    1.3/70.9 MB 6.3 MB/s eta
0:00:12
                                    1.7/70.9 MB 6.4 MB/s eta
0:00:11
       -                            2.2/70.9 MB 7.3 MB/s eta
0:00:10
       -                            2.7/70.9 MB 7.8 MB/s eta
0:00:09
       -                            2.8/70.9 MB 7.8 MB/s eta
0:00:09
```

```python
from xgboost import XGBClassifier
#create a copy of the cleaned dataset

#define dependent variable and independent variable

#Split the dataset into the Training and the Test set. Set the test set to 0.3
# Apply Xgboost
model = XGBClassifier()

#fit model
model.fit(X_train, y_train)
# make predictions for test data
y_pred_xgb = model.predict(X_test)

# evaluate predictions
predictions = [round(value) for value in y_pred_xgb]
accuracy = accuracy_score(y_test, predictions)
print("Accuracy: %.2f%%" % (accuracy * 100.0))
```

```
Accuracy: 79.78%
```

# Out of all 3 approaches, which gives you the best result?

```python
#Random Forest Classifier
```