# Question1

```python
import time
from sklearn.svm import SVC
from sklearn.linear_model import LogisticRegression
from sklearn.ensemble import AdaBoostClassifier
from sklearn.ensemble import RandomForestClassifier
from sklearn.tree import DecisionTreeClassifier
from sklearn.neural_network import MLPClassifier
from sklearn.model_selection import train_test_split
from sklearn.preprocessing import StandardScaler
from sklearn.datasets import make_classification
```

```python
import numpy as np
import matplotlib.pyplot as plt
from matplotlib.colors import ListedColormap
import warnings

warnings.simplefilter(action='ignore', category=FutureWarning)
warnings.filterwarnings("ignore")
```

```python
def create_dataset(n=1250, nf=2, nr=0, ni=2, random_state=125):
    '''
    generate a new dataset with
    n: total number of samples
    nf: number of features
    nr: number of redundant features (these are linear combinatins of
    informative features)
    ni: number of informative features (ni + nr = nf must hold)
    random_state: set for reproducibility
    '''
    X, y = make_classification(n_samples=n,
                               n_features=nf,
                               n_redundant=nr,
                               n_informative=ni,
                               random_state=random_state,
                               n_clusters_per_class=2)
    rng = np.random.RandomState(2)
    X += 3 * rng.uniform(size=X.shape)
    X = StandardScaler().fit_transform(X)
    return X, y
```

```python
def plotter(classifier, X, X_test, y_test, title, ax=None):
    # plot decision boundary for given classifier
```

```
 3        plot_step = 0.02
 4        x_min, x_max = X[:, 0].min() - 1, X[:, 0].max() + 1
 5        y_min, y_max = X[:, 1].min() - 1, X[:, 1].max() + 1
 6        xx, yy = np.meshgrid(np.arange(x_min, x_max, plot_step),
 7                             np.arange(y_min, y_max, plot_step))
 8        Z = classifier.predict(np.c_[xx.ravel(), yy.ravel()])
 9        Z = Z.reshape(xx.shape)
10        if ax:
11            ax.contourf(xx, yy, Z, cmap=plt.cm.Paired)
12            ax.scatter(X_test[:, 0], X_test[:, 1], c=y_test)
13            ax.set_title(title)
14        else:
15            plt.contourf(xx, yy, Z, cmap=plt.cm.Paired)
16            plt.scatter(X_test[:, 0], X_test[:, 1], c=y_test)
17            plt.title(title)
```
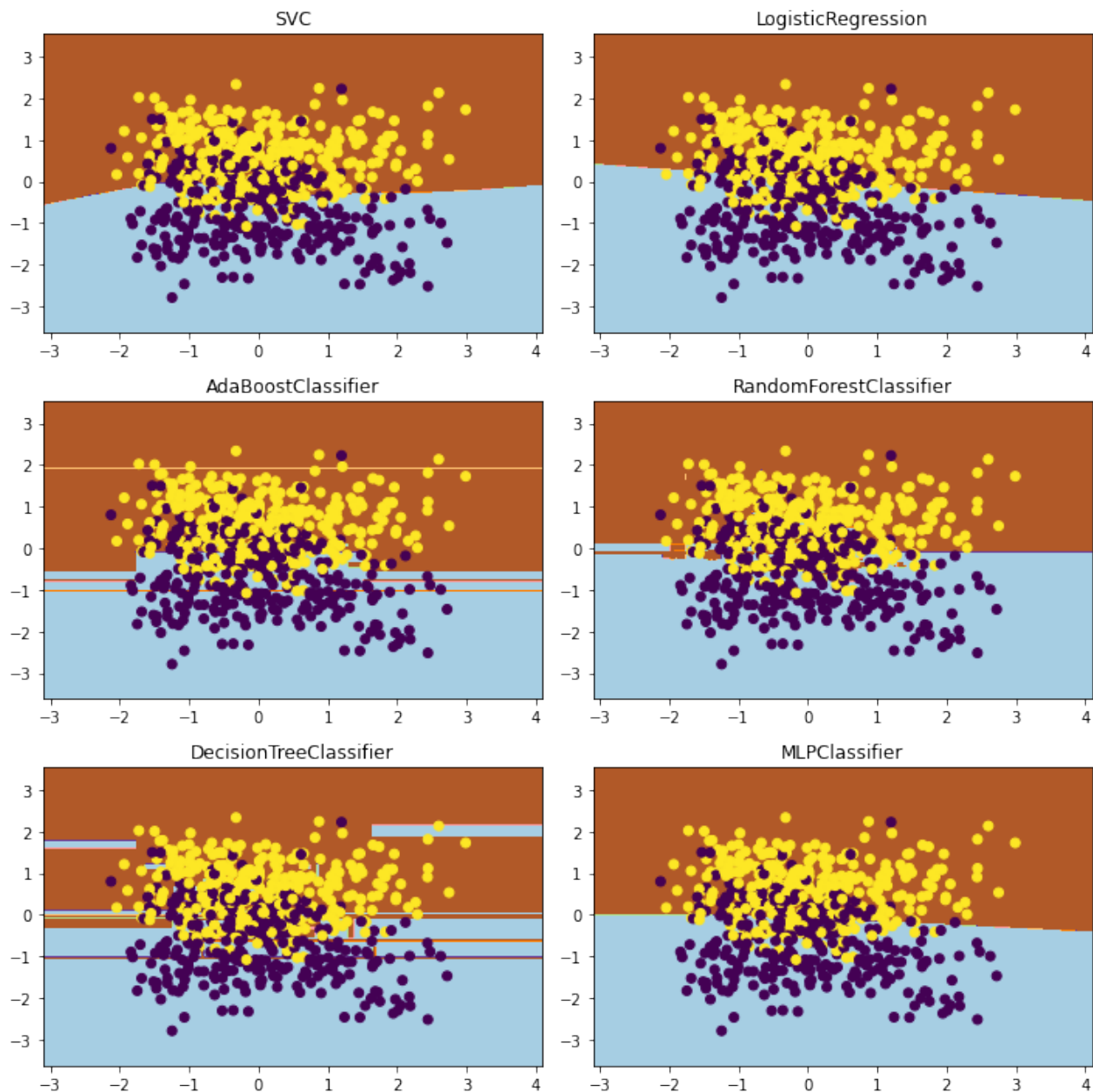
**(a)**

```
 1  # data set
 2  data, target = create_dataset()
 3
 4  # split
 5  X_train, X_test, Y_train, Y_test = train_test_split(data, target,
    train_size=0.5, test_size=0.5, random_state=15)
 6
 7  # classifier
 8  classified_list = []
 9
10  SVCModel = SVC()
11  SVCModel.fit(X_train, Y_train)
12  classified_list.append(SVCModel)
13
14  LogisticRegressionModel = LogisticRegression()
15  LogisticRegressionModel.fit(X_train, Y_train)
16  classified_list.append(LogisticRegressionModel)
17
18  AdaBoostClassifierModel = AdaBoostClassifier()
19  AdaBoostClassifierModel.fit(X_train, Y_train)
20  classified_list.append(AdaBoostClassifierModel)
21
22  RandomForestClassifierModel = RandomForestClassifier()
23  RandomForestClassifierModel.fit(X_train, Y_train)
24  classified_list.append(RandomForestClassifierModel)
25
26  DecisionTreeClassifierModel = DecisionTreeClassifier()
27  DecisionTreeClassifierModel.fit(X_train, Y_train)
28  classified_list.append(DecisionTreeClassifierModel)
29
```

```python
30  MLPClassifierModel = MLPClassifier()
31  MLPClassifierModel.fit(X_train, Y_train)
32  classified_list.append(MLPClassifierModel)
33
34  classified_name_list = ['SVC', 'LogisticRegression', 'AdaBoostClassifier',
    'RandomForestClassifier',
35                          'DecisionTreeClassifier', 'MLPClassifier']
36
37  fig, ax = plt.subplots(3, 2, figsize=(10, 10))
38  for i, ax in enumerate(ax.flat):
39      plotter(classifier=classified_list[i], X=X_train, X_test=X_test,
    y_test=Y_test, title=classified_name_list[i],
40              ax=ax)
41  plt.tight_layout()
42  plt.show()
```

**(b)**

```
1   # data set
2   data, target = create_dataset()
3
4   X_train, X_test, Y_train, Y_test = train_test_split(data, target,
    train_size=0.8, test_size=0.2, random_state=45)
5   x_y_set = np.hstack((X_train, Y_train.reshape(1000, 1)))
6
7   size_list = [50, 100, 200, 300, 400, 500, 600, 700, 800, 900, 1000]
8
9   # np.random.shuffle(x_y_set)
10
11  # Decision Tree
```
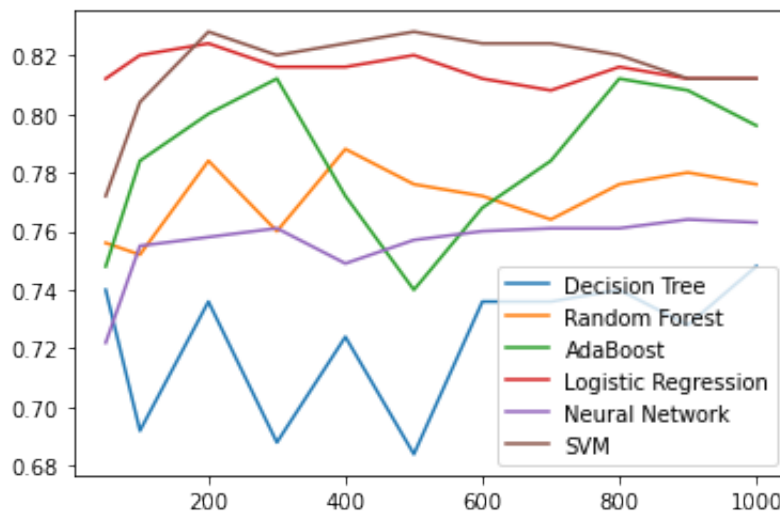
```python
12   DecisionTreeClassifierModel = DecisionTreeClassifier()
13   decision_tree_list = []
14   for i in range(len(size_list)):
15       np.random.shuffle(x_y_set)
16       DecisionTreeClassifierModel.fit(x_y_set[:size_list[i], :2],
     x_y_set[:size_list[i], 2])
17       decision_tree_list.append(DecisionTreeClassifierModel.score(X_test,
     Y_test))
18
19   # Random Forest
20   RandomForestClassifierModel = RandomForestClassifier()
21   random_forest_list = []
22   for i in range(len(size_list)):
23       np.random.shuffle(x_y_set)
24       RandomForestClassifierModel.fit(x_y_set[:size_list[i], :2],
     x_y_set[:size_list[i], 2])
25       random_forest_list.append(RandomForestClassifierModel.score(X_test,
     Y_test))
26
27   # AdaBoost
28   AdaBoostClassifierModel = AdaBoostClassifier()
29   ada_boost_list = []
30   for i in range(len(size_list)):
31       np.random.shuffle(x_y_set)
32       AdaBoostClassifierModel.fit(x_y_set[:size_list[i], :2],
     x_y_set[:size_list[i], 2])
33       ada_boost_list.append(AdaBoostClassifierModel.score(X_test, Y_test))
34
35   # LogisticRegression
36   LogisticRegressionModel = LogisticRegression()
37   logistic_regression_list = []
38   for i in range(len(size_list)):
39       np.random.shuffle(x_y_set)
40       LogisticRegressionModel.fit(x_y_set[:size_list[i], :2],
     x_y_set[:size_list[i], 2])
41       logistic_regression_list.append(LogisticRegressionModel.score(X_test,
     Y_test))
42
43   # MLPClassifier
44   MLPClassifierModel = MLPClassifier()
45   neural_network_list = []
46   for i in range(len(size_list)):
47       np.random.shuffle(x_y_set)
48       MLPClassifierModel.fit(x_y_set[:size_list[i], :2],
     x_y_set[:size_list[i], 2])
49       neural_network_list.append(MLPClassifierModel.score(X_train, Y_train))
50
51   # SVC
52   SVCModel = SVC()
```

```
53    svc_list = []
54    for i in range(len(size_list)):
55        # np.random.shuffle(x_y_set)
56        SVCModel.fit(x_y_set[:size_list[i], :2], x_y_set[:size_list[i], 2])
57        svc_list.append(SVCModel.score(X_test, Y_test))
58
59    plt.plot(size_list, decision_tree_list, label='Decision Tree')
60    plt.plot(size_list, random_forest_list, label='Random Forest')
61    plt.plot(size_list, ada_boost_list, label='AdaBoost')
62    plt.plot(size_list, logistic_regression_list, label='Logistic Regression')
63    plt.plot(size_list, neural_network_list, label='Neural Network')
64    plt.plot(size_list, svc_list, label='SVM')
65    plt.legend()
66    plt.show()
```



From the graph, i will choose Logistic Regression model because it has higher accuracy and more stable than others.

## (c)

```
1    # data set
2    data, target = create_dataset()
3
4    # split
5    # X_train, X_test = train_test_split(data, train_size=0.8, test_size=0.2,
     random_state=45)
6    # Y_train, Y_test = train_test_split(target, train_size=0.8,
     test_size=0.2, random_state=45)
7
8    X_train, X_test, Y_train, Y_test = train_test_split(data, target,
     train_size=0.8, test_size=0.2, random_state=45)
9    # x_y_set = np.array()
```
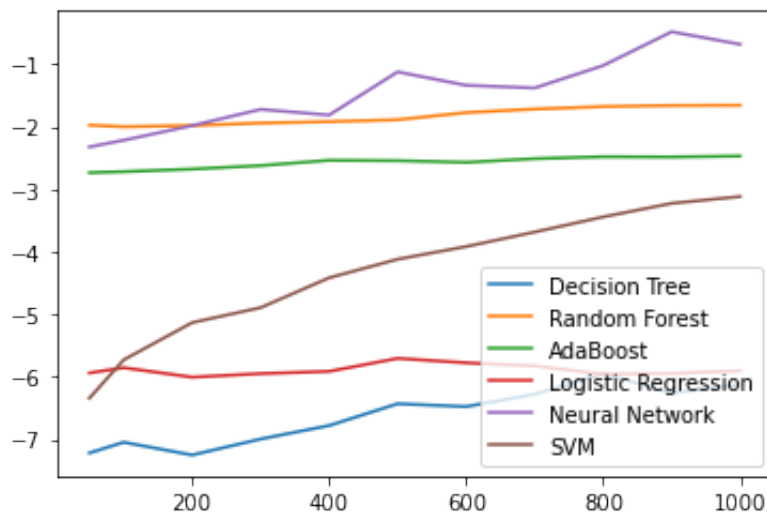
```python
10
11  size_list = [50, 100, 200, 300, 400, 500, 600, 700, 800, 900, 1000]
12
13  # Decision Tree
14  DecisionTreeClassifierModel = DecisionTreeClassifier()
15  decision_tree_list = []
16  dt_time_list = []
17  for i in range(len(size_list)):
18      start = time.time()
19      DecisionTreeClassifierModel.fit(X_train[:size_list[i]],
    Y_train[:size_list[i]])
20      decision_tree_list.append(DecisionTreeClassifierModel.score(X_test,
    Y_test))
21      end = time.time()
22      dt_time_list.append(np.log(end - start))
23
24  # Random Forest
25  RandomForestClassifierModel = RandomForestClassifier()
26  random_forest_list = []
27  rf_time_list = []
28  for i in range(len(size_list)):
29      start = time.time()
30      RandomForestClassifierModel.fit(X_train[:size_list[i]],
    Y_train[:size_list[i]])
31      random_forest_list.append(RandomForestClassifierModel.score(X_test,
    Y_test))
32      end = time.time()
33      rf_time_list.append(np.log(end - start))
34
35  # AdaBoost
36  AdaBoostClassifierModel = AdaBoostClassifier()
37  ada_boost_list = []
38  ab_time_list = []
39  for i in range(len(size_list)):
40      start = time.time()
41      AdaBoostClassifierModel.fit(X_train[:size_list[i]],
    Y_train[:size_list[i]])
42      ada_boost_list.append(AdaBoostClassifierModel.score(X_test, Y_test))
43      end = time.time()
44      ab_time_list.append(np.log(end - start))
45
46  # LogisticRegression
47  LogisticRegressionModel = LogisticRegression()
48  logistic_regression_list = []
49  lr_time_list = []
50  for i in range(len(size_list)):
51      start = time.time()
52      LogisticRegressionModel.fit(X_train[:size_list[i]],
    Y_train[:size_list[i]])
```

```python
53          logistic_regression_list.append(LogisticRegressionModel.score(X_test,
     Y_test))
54      end = time.time()
55      lr_time_list.append(np.log(end - start))
56
57  # MLPClassifier
58  MLPClassifierModel = MLPClassifier()
59  neural_network_list = []
60  nn_time_list = []
61  for i in range(len(size_list)):
62      start = time.time()
63      MLPClassifierModel.fit(X_train[:size_list[i]], Y_train[:size_list[i]])
64      neural_network_list.append(MLPClassifierModel.score(X_train, Y_train))
65      end = time.time()
66      nn_time_list.append(np.log(end - start))
67
68  # SVC
69  SVCModel = SVC()
70  svc_list = []
71  svc_time_list = []
72  for i in range(len(size_list)):
73      start = time.time()
74      SVCModel.fit(X_train[:size_list[i]], Y_train[:size_list[i]])
75      svc_list.append(SVCModel.score(X_test, Y_test))
76      end = time.time()
77      svc_time_list.append(np.log(end - start))
78
79  plt.plot(size_list, dt_time_list, label='Decision Tree')
80  plt.plot(size_list, rf_time_list, label='Random Forest')
81  plt.plot(size_list, ab_time_list, label='AdaBoost')
82  plt.plot(size_list, lr_time_list, label='Logistic Regression')
83  plt.plot(size_list, nn_time_list, label='Neural Network')
84  plt.plot(size_list, svc_time_list, label='SVM')
85  plt.legend()
86  plt.show()
```

Random Forest and AdaBoost is stable and fast while incresing the number of data, SVM model's time get more while data amount incresing, Decision Tree and Logistic Regression is slow and hradly affected by amount of data.

## (d)

```
1   data, target = create_dataset(n=2000, nf=20, nr=12, ni=8, random_state=25)
2   X_train, X_test, Y_train, Y_test = train_test_split(data, target,
    train_size=0.5, test_size=0.5, random_state=15)
3   DecisionTreeClassifierModel = DecisionTreeClassifier()
4   DecisionTreeClassifierModel.fit(X_train, Y_train)
5
6   train_accuracy = DecisionTreeClassifierModel.score(X_train, Y_train)
7   test_accuracy = DecisionTreeClassifierModel.score(X_test, Y_test)
8
9   print('train_accuracy: ', train_accuracy)
10  print('test_accuracy: ', test_accuracy)
11
```

```
1   train_accuracy:   1.0
2   test_accuracy:   0.832
```
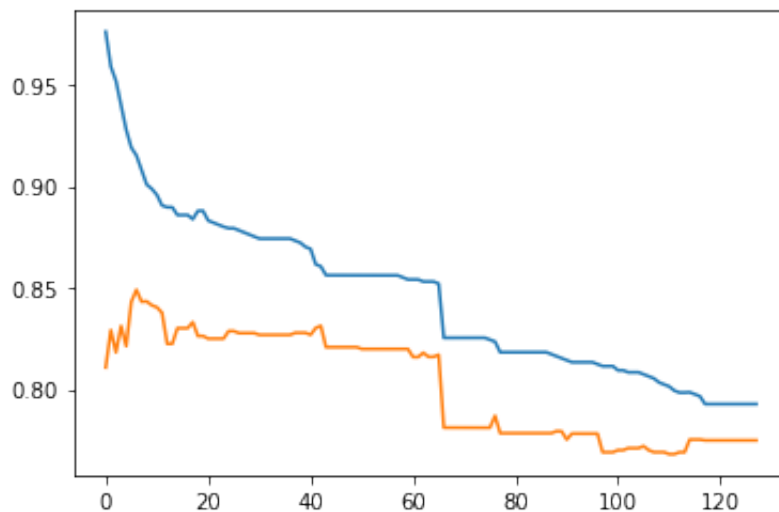
## (e)

```
1   from sklearn.metrics import roc_curve
2   from sklearn.metrics import auc
3
4   data, target = create_dataset(n=2000, nf=20, nr=12, ni=8, random_state=25)
5   X_train, X_test, Y_train, Y_test = train_test_split(data, target,
    train_size=0.5, test_size=0.5, random_state=15)
6   auc_train_list = []
```

```
 7   auc_test_list = []
 8   for i in range(2, 130):
 9       DecisionTreeClassifierModel =
     DecisionTreeClassifier(min_samples_leaf=i)
10       DecisionTreeClassifierModel.fit(X_train, Y_train)
11       train_predict = DecisionTreeClassifierModel.predict(X_train)
12       fpr_train, tpr_train, _ = roc_curve(Y_train, train_predict)
13       auc_train_list.append(auc(fpr_train, tpr_train))
14
15       test_predict = DecisionTreeClassifierModel.predict(X_test)
16       fpr_test, tpr_test, _ = roc_curve(Y_test, test_predict)
17       auc_test_list.append(auc(fpr_test, tpr_test))
18
19   plt.plot(auc_train_list)
20   plt.plot(auc_test_list)
21   plt.show()
```



## (f)

```
 1   data, target = create_dataset(n=2000, nf=20, nr=12, ni=8, random_state=25)
 2   X_train, X_test, Y_train, Y_test = train_test_split(data, target,
     train_size=0.5, test_size=0.5, random_state=15)
 3   X_train_split = np.split(X_train, 10)
 4   X_test_split = np.split(X_test, 10)
 5   Y_train_split = np.split(Y_train, 10)
 6   Y_test_split = np.split(Y_test, 10)
 7   score = []
 8   x_axis = [i for i in range(2, 96)]
 9   test_accuracy = 0
10   train_accuracy = 0
11   max_auc = 0
12   best_k = 0
```
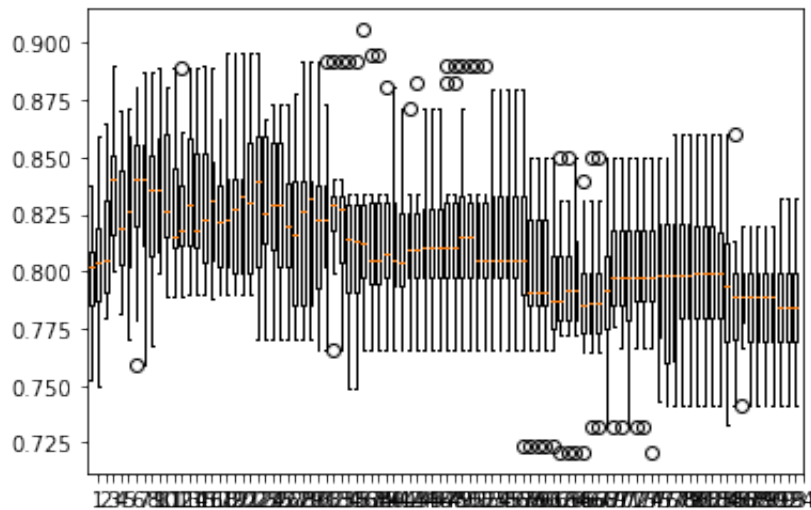
```python
for i in range(2, 96):
    k = 0
    temp = []
    for j in range(10):
        new_x = []
        new_y = []
        DecisionTreeClassifierModel = DecisionTreeClassifier(min_samples_leaf=i)
        for m in range(10):
            if m == k:
                continue
            new_x.append(X_train_split[m])
            new_y.append(Y_train_split[m])
        new_x = np.concatenate(new_x, axis=0)
        new_y = np.concatenate(new_y, axis=0)
        DecisionTreeClassifierModel.fit(new_x, new_y)
        train_predict = DecisionTreeClassifierModel.predict(X_train_split[k])
        ftr, tpr, _ = roc_curve(Y_train_split[k], train_predict)
        temp_auc = auc(ftr, tpr)
        if temp_auc > max_auc:
            max_auc = temp_auc
            train_accuracy = DecisionTreeClassifierModel.score(X_train, Y_train)
            test_accuracy = DecisionTreeClassifierModel.score(X_test, Y_test)
            best_k = i
        temp.append(temp_auc)
        k += 1
    score.append(temp)
    k = 0

print('min_samples_leaf: ', best_k)
print('train_accuracy: ', train_accuracy)
print('test_accuracy: ', test_accuracy)
plt.boxplot(score)
plt.xticks(np.arange(2, 96, 1))
plt.show()
```

```
min_samples_leaf:  38
train_accuracy:  0.862
test_accuracy:  0.831
```

## (g)

```
1   from sklearn.model_selection import GridSearchCV
2
3   data, target = create_dataset(n=2000, nf=20, nr=12, ni=8, random_state=25)
4   X_train, X_test, Y_train, Y_test = train_test_split(data, target,
    train_size=0.5, test_size=0.5, random_state=15)
5   X_train_split = np.split(X_train, 10)
6   X_test_split = np.split(X_test, 10)
7   Y_train_split = np.split(Y_train, 10)
8   Y_test_split = np.split(Y_test, 10)
9
10  para = {'min_samples_leaf': [i for i in range(2, 96)]}
11  clf = GridSearchCV(param_grid=para, scoring='roc_auc', cv=10,
    estimator=DecisionTreeClassifier())
12  clf_res = clf.fit(X_train, Y_train)
13  train_accuracy = clf_res.score(X_train, Y_train)
14  test_accuracy = clf_res.score(X_test, Y_test)
15  # print(clf_res.best_score_)
16  print(clf_res.best_params_)
17  print('train_accuracy: ', train_accuracy)
18  print('test_accuracy: ', test_accuracy)
```

```
1   {'min_samples_leaf': 27}
2   train_accuracy:  0.953077307730773
3   test_accuracy:  0.9007445459310025
```

Becuase of different data processing methods, gridsearch is random while in (f) we use the k-fold.