# 课后作业：决策树(Decision Tree)与随机森林(Random Forests)

作者：欧新宇 (Xinyu OU)

本文档所展示的测试结果，均运行于：Intel Core i7-7700K CPU 4.2GHz

**【作业提交】**

将分类结果保存到文本文档进行提交(写上每一题的题号和题目，然后再贴答案)，同时提交源代码。

1. 测试结果命名为: ex06-结果-你的学号-你的姓名.txt
2. 输出图片命名为: ex06-性能对比图-你的学号-你的姓名.png (.jpg)
3. 源代码命名为: ex06-01-你的学号-你的姓名.py, ex06-02-你的学号-你的姓名.py, ex06-03-你的学号-你的姓名.py

*结果文件，要求每小题标注题号，两题之间要求空一行*

---

要求在 "糖尿病预测" 数据集上分别使用决策树与随机森林完成以下任务，要求如下：

1. 要求训练集和测试集的分割比例为75%:25%
2. 使用**决策树**模型输出树的深度分别为3和5的得分，要求同时输出训练集和测试集上的评分结果。 (ex06-01)
3. 使用**随机森林**模型输出森林中树的个数分别为4和6的得分，随机数种子=8，要求同时输出训练集和测试集上的评分结果。 (ex06-02)
4. 同时使用**决策树**(树深度={1:20})和**随机森林**(树的棵树={1:20})进行建模，并输出性能对比图。 (ex06-03, ex06-性能对比图)

- **决策树**

```python
# 加载 pandas库，并使用read_csv()函数读取糖尿病预测数据集diabetes
import pandas as pd
from sklearn import tree
from sklearn.model_selection import train_test_split

data = pd.read_csv('../Datasets/diabetes.csv')


# 将数据中的特征和标签进行分离，其中第0位位索引号，第1-8位位特征，第9位为标签
X = data.iloc[:, 0:8]
y = data.iloc[:, 8]

# 以 70%:30%的比例对训练集和测试集进行拆分
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size = 0.25)

dt3 = tree.DecisionTreeClassifier(max_depth = 3)
dt5 = tree.DecisionTreeClassifier(max_depth = 5)
dt3.fit(X_train, y_train)
dt5.fit(X_train, y_train)

print("max_depth=3，训练集评分:{0:.3f}；测试集评分:{1:.3f}".format(dt3.score(X_train, y_train), dt3.score(X_test, y_test)))
```

```
22  print("max_depth=5，训练集评分:{0:.3f}；测试集评分:
    {1:.3f}".format(dt5.score(X_train, y_train), dt5.score(X_test, y_test)))
23
24
```

```
1  max_depth=3，训练集评分:0.783；测试集评分:0.698
2  max_depth=5，训练集评分:0.844；测试集评分:0.698
```

- **随机森林**

```
1   # 加载 pandas库，并使用read_csv()函数读取糖尿病预测数据集diabetes
2   import pandas as pd
3   from sklearn.ensemble import RandomForestClassifier
4   from sklearn.model_selection import train_test_split
5
6   data = pd.read_csv('../Datasets/diabetes.csv')
7
8   # 将数据中的特征和标签进行分离，其中第0位位索引号，第1-8位位特征，第9位为标签
9   X = data.iloc[:, 0:8]
10  y = data.iloc[:, 8]
11
12  # 以 70%:30%的比例对训练集和测试集进行拆分
13  X_train, X_test, y_train, y_test = train_test_split(X, y, test_size = 0.25)
14
15
16  rf3 = RandomForestClassifier(n_estimators = 3, random_state = 8, n_jobs =
    -1)
17  rf5 = RandomForestClassifier(n_estimators = 5, random_state = 8, n_jobs =
    -1)
18  rf3.fit(X_train, y_train)
19  rf5.fit(X_train, y_train)
20
21  print("n_estimators=3，训练集评分:{0:.3f}；测试集评分:
    {1:.3f}".format(rf3.score(X_train, y_train), rf3.score(X_test, y_test)))
22  print("n_estimators=5，训练集评分:{0:.3f}；测试集评分:
    {1:.3f}".format(rf5.score(X_train, y_train), rf5.score(X_test, y_test)))
23
24
```

```
1  n_estimators=3，训练集评分:0.941；测试集评分:0.714
2  n_estimators=5，训练集评分:0.970；测试集评分:0.698
```

- **性能对比**

```
1   # 加载 pandas库，并使用read_csv()函数读取糖尿病预测数据集diabetes
2   import numpy as np
3   import pandas as pd
4   from sklearn import tree
5   from sklearn.ensemble import RandomForestClassifier
6   from sklearn.model_selection import train_test_split
7   import matplotlib.pyplot as plt
8   import os
9
10  data = pd.read_csv(os.path.join(os.getcwd(), '..', 'datasets',
    'diabetes.csv'))
```

```python
11
12  # 将数据中的特征和标签进行分离，其中第0位位索引号，第1-8位位特征，第9位为标签
13  X = data.iloc[:, 0:8]
14  y = data.iloc[:, 8]
15
16  # 以 70%:30%的比例对训练集和测试集进行拆分
17  X_train, X_test, y_train, y_test = train_test_split(X, y, test_size = 0.25)
18
19  n = 40
20  scores = np.zeros([4, n]) #第1-4列分别为:
    score_train_dt,score_test_dt,score_train_rf,score_test_rf
21  num = np.arange(0, n)
22
23  for i in num:
24      n = i + 1
25
26      # 利用当行刷新方法显示正在计算的模型
27      print("\r 正在计算第{}/{}个模型，请稍等...".format(n, num.shape[0]),
    end="")
28
29      dt = tree.DecisionTreeClassifier(max_depth = n)
30      dt.fit(X_train, y_train)
31
32      rf = RandomForestClassifier(n_estimators = n, random_state = 8, n_jobs
    = -1)
33      rf.fit(X_train, y_train)
34
35      scores[0, i] = dt.score(X_train, y_train)
36      scores[1, i] = dt.score(X_test, y_test)
37      scores[2, i] = rf.score(X_train, y_train)
38      scores[3, i] = rf.score(X_test, y_test)
39
40  #      print("随机森林的评分:{}.".format(rf.score(X_test, y_test)))
41
42      if i == num.shape[0] - 1:
43          print("计算完毕! ")
44
45
46
47  plt.figure(dpi=100)
48  plt.plot(num, scores[0,:], label="DecisionTree_Train")
49  plt.plot(num, scores[1,:], label="DecisionTree_Test")
50  plt.plot(num, scores[2,:], label="ForestClassifier_Train")
51  plt.plot(num, scores[3,:], label="ForestClassifier_Test")
52
53  plt.legend(loc='upper right')
54  plt.savefig('results/Ch06Hw01DecisionTree.png', dpi=150)
55  plt.show()
```

```
1  正在计算第40/40个模型，请稍等...计算完毕!
```