

## 1、处理数据

- 删去某些特征缺失数据的数据
- 用随机森林补全那些缺失数据多的数据 ( 图片：处理数据代码.png)

```

1: 1 # 处理数据
2   import numpy as np
3   import pandas as pd
4
5   # 读取文件
6   data_set = pd.read_csv('train.csv', header=None)
7   # 将? 用NULL 来替换
8   new_data_set = data_set.replace(['?'], [None])
9   new_data_set.to_csv('new_data_set1.csv')
10  """
11  ## 用随机森林补全数据
12  * 特征数据完整的特征 0 2 6 8 13
13  * 因为有 7000 多条数据, 对于特征缺少比较小的部分的数据进行删除
14  * * 1 3 4 9 11
15  * 对于 5 7 10 12 缺失较多的特征用随机森林来补全
16  """
17  # 删数据
18  new_data_set = pd.read_csv('new_data_set.csv')
19  new_data_set.dropna(subset=['1', '3', '4', '9', '11'], inplace=True)
20
21  # 利用随机森林补全 12 列
22  from sklearn.ensemble import RandomForestRegressor
23
24  def set_missing_ages(df):
25      # 把已有数值特征取出来丢进Random Forest Regressor (随机森林回归) 中
26      all_df = df[['0', '1', '2', '3', '4', '6', '8', '9', '11', '12', '13']]
27      know_goal = all_df[all_df['12'].notnull()]
28      unknow_goal = all_df[all_df['12'].isnull()]
29
30      # 构建随机森林模型
31      y = know_goal['12'].values.tolist()
32      know_goal = know_goal.drop('12', 1)
33      x = know_goal.values
34
35      rfr = RandomForestRegressor(random_state=0, n_estimators=2000, n_jobs=-1)
36      rfr.fit(x, y)
37
38      un_x = unknow_goal.drop('12', 1).values
39      pred = rfr.predict(un_x)
40
41      for i in range(len(pred)):
42          pred[i] = int(pred[i])
43
44      df.loc[(df['12'].isnull()), '12'] = pred
45
46      return df
47  new_data_set = set_missing_ages(new_data_set)
48  # 10
49  def set_missing_ages(df):
50      # 把已有数值特征取出来丢进Random Forest Regressor (随机森林回归) 中
51      all_df = df[['0', '1', '2', '3', '4', '6', '8', '9', '10', '11', '12', '13']]
52      know_goal = all_df[all_df['10'].notnull()]
53      unknow_goal = all_df[all_df['10'].isnull()]
54
55      # 构建随机森林模型

```

```

56     y = know_goal['10'].values.tolist()
57     know_goal = know_goal.drop('10',1)
58     x = know_goal.values
59
60     rfr = RandomForestRegressor(random_state=0,n_estimators=2000,n_jobs=-1)
61     rfr.fit(x,y)
62
63     un_x = unknow_goal.drop('10',1).values
64     pred = rfr.predict(un_x)
65
66     for i in range(len(pred)):
67         pred[i] = int(pred[i])
68
69     df.loc[(df['10'].isnull()),'10'] = pred
70
71     return df
72 new_data_set = set_missing_ages(new_data_set)
73
74 # 7
75 def set_missing_ages(df):
76     # 把已有数值特征取出来丢进Random Forest Regressor (随机森林回归) 中
77     all_df = df[['0', '1', '2', '3', '4', '6', '7', '8', '9', '10', '11', '12', '13']]
78     know_goal = all_df[all_df['7'].notnull()]
79     unknow_goal = all_df[all_df['7'].isnull()]
80
81     # 构建随机森林模型
82     y = know_goal['7'].values.tolist()
83     know_goal = know_goal.drop('7',1)
84     x = know_goal.values
85
86     rfr = RandomForestRegressor(random_state=0,n_estimators=2000,n_jobs=-1)
87     rfr.fit(x,y)
88
89     un_x = unknow_goal.drop('7',1).values
90     pred = rfr.predict(un_x)
91
92     for i in range(len(pred)):
93         pred[i] = int(pred[i])
94
95     df.loc[(df['7'].isnull()),'7'] = pred
96
97     return df
98 new_data_set = set_missing_ages(new_data_set)
99
100
101 # 5
102 def set_missing_ages(df):
103     # 把已有数值特征取出来丢进Random Forest Regressor (随机森林回归) 中
104     all_df = df[['0', '1', '2', '3', '4', '5', '6', '8', '9', '10', '11', '12', '13']]
105     know_goal = all_df[all_df['5'].notnull()]
106     unknow_goal = all_df[all_df['5'].isnull()]
107
108     # 构建随机森林模型
109     y = know_goal['5'].values.tolist()
110     know_goal = know_goal.drop('5',1)
111     x = know_goal.values
112
113     rfr = RandomForestRegressor(random_state=0,n_estimators=2000,n_jobs=-1)
114     rfr.fit(x,y)
115
116     un_x = unknow_goal.drop('5',1).values

```

```

117     pred = rfr.predict(un_x)
118
119     for i in range(len(pred)):
120         pred[i] = int(pred[i])
121
122     df.loc[(df['5'].isnull()), '5'] = pred
123
124     return df
125 new_data_set = set_missing_ages(new_data_set)
126 new_data_set.info()
127
128 # 保存处理好的数据
129 new_data_set.to_csv('complete_data_set1.csv')
130

```

executed in 36.8s, finished 11:38:44 2020-10-23

```

<class 'pandas.core.frame.DataFrame'>
Int64Index: 6694 entries, 0 to 7192
Data columns (total 15 columns):
#   Column      Non-Null Count  Dtype
---  ---
0   Unnamed: 0    6694 non-null   int64
1   0             6694 non-null   int64
2   1             6694 non-null   float64
3   2             6694 non-null   int64
4   3             6694 non-null   float64
5   4             6694 non-null   float64
6   5             6694 non-null   float64
7   6             6694 non-null   int64
8   7             6694 non-null   float64
9   8             6694 non-null   int64
10  9             6694 non-null   float64
11  10            6694 non-null   float64
12  11            6694 non-null   float64
13  12            6694 non-null   float64
14  13            6694 non-null   int64
dtypes: float64(9), int64(6)
memory usage: 836.8 KB

```

## 2、挑选模型

- 为了保证模型的准确率，所以使用5折交叉验证的方法来选取最佳模型

### KNN模型

- 代码：

```
In [ ]: 1 from sklearn.model_selection import train_test_split
2 from sklearn.ensemble import RandomForestClassifier
3 from sklearn.model_selection import GridSearchCV
4 from sklearn import metrics
5 # 分层划分数据
6 from sklearn.model_selection import StratifiedKFold
7 # 获取数据
8 datas = pd.read_csv("complete_data_set.csv")
9 datas = datas.drop('Unnamed: 0',1).values
10 X = datas[:, :-1]
11 Y = datas[:, -1]
12 # 分层划分数据，采用五折交叉验证
13 sfolder = StratifiedKFold(n_splits=5, shuffle=True, random_state=13)
14
15 # 利用cross_val_score自动获得结果
16 from sklearn.neighbors import KNeighborsClassifier
17 from sklearn.model_selection import cross_val_score
18 from sklearn.preprocessing import MinMaxScaler
19 min_max = MinMaxScaler()
20 X = min_max.fit_transform(X)
21 print(X[0])
22
23 max_ = 0.0
24 string = None
25
26 k_list = [i for i in range(1, 480, 2)]
27 for k in k_list:
28     clf=KNeighborsClassifier(n_neighbors=k)
29     scores=cross_val_score(clf, X, Y, cv=sfolder)
30     print("k = "+str(k), end=" ")
31     print(scores, end=" ")
32     avg_s = np.mean(np.array(scores))
33     print(avg_s)
34     if max_ < avg_s:
35         max_ = avg_s
36         string = ("{} {} {} ->{}".format(c, depth, feat, avg_s))
37
38 print("最好： "+string)
```

- 运行最好效果参数
  - 最好：47 ->0.33179033948767067

## 随机森林

- 代码

"""key

```
# 随机森林
# 利用cross_val_score自动获得结果
import numpy as np
import pandas as pd
from sklearn.model_selection import cross_val_score
# 随机森林
from sklearn.ensemble import RandomForestClassifier

# 获取数据
datas = pd.read_csv("complete_data_set.csv")
datas = datas.drop('Unnamed: 0',1).values
X = datas[:, :-1]
Y = datas[:, -1]
# 分层划分数据，采用五折交叉验证
```

```

sfolder = StratifiedKFold(n_splits=5,shuffle=True, random_state=13)

criterion_list = ['gini','entropy']
max_depth_list = [i for i in range(4,16)]
max_features_list = ['auto','sqrt','log2']

max_ = 0.0
string = None
for c in criterion_list:
    for depth in max_depth_list:
        for feat in max_features_list:
            clf =
RandomForestClassifier(criterion=c,max_depth=depth,max_features=feat,n_estimators=
120)

            scores=cross_val_score(clf, X, Y, cv=sfolder)
            avg_s = np.mean(np.array(scores))
            print("{} {} {}  ->{}".format(c,depth,feat,avg_s))
            if max_ < avg_s:
                max_ = avg_s
                string = ("{} {} {}  ->{}".format(c,depth,feat,avg_s))

print("最好： "+string)

```

'''

- 运行最好效果参数
- ◦ gini 8 log2 ->0.3612178510389142

## SVM支持向量机

- 代码

'''

```

import numpy as np
from sklearn.pipeline import make_pipeline
from sklearn.preprocessing import StandardScaler
from sklearn.svm import SVC
from sklearn.model_selection import cross_val_score

# 获取数据
datas = pd.read_csv("/kaggle/input/123345/complete_data_set.csv")
datas = datas.drop('Unnamed: 0',1).values
X = datas[:, :-1]
Y = datas[:, -1]
# 分层划分数据，采用五折家产验证
sfolder = StratifiedKFold(n_splits=5,shuffle=True, random_state=13)
# 对数据进行归一化
from sklearn.preprocessing import MinMaxScaler
min_max = MinMaxScaler()

```

```

X = min_max.fit_transform(X)

C_list = [1e-3, 1e-2, 1e-1, 1, 10, 100, 1000]
kernel_list = ['poly', 'rbf', 'sigmoid', 'precomputed']
gamma_list = ['scale', 'auto']
class_weight_list = ['balanced']

max_ = 0.0
string = None

for k in kernel_list:
    for gamma in gamma_list:
        for c in C_list:
            for class_weight in class_weight_list:
                svc = SVC(kernel=k, gamma=gamma, C=c, class_weight=class_weight)
                clf = make_pipeline(StandardScaler(), svc)
                scores = cross_val_score(clf, X, Y, cv=sfolds)
                avg_s = np.mean(np.array(scores))
                print("{} {} {} {} -> {}".format(k, gamma, c, class_weight, avg_s))
                if max_ < avg_s:
                    max_ = avg_s
                    string = ("{} {} {} {} -> {}".format(k, gamma, c, class_weight, avg_s))

print(string)

```

'''

- 运行最好效果参数
- ◦ rbg auto 1 balanced -> 0.3462178510389142(应该是这个参数，当时跑了好几个小时得出来的)

## 神经网络

- 代码

'''key

```

# 利用cross_val_score自动获得结果
import numpy as np
import pandas as pd
from sklearn.model_selection import cross_val_score, StratifiedKFold
# 获取数据
from sklearn.neural_network import MLPClassifier
from sklearn.preprocessing import MinMaxScaler
min_max = MinMaxScaler()

import warnings

warnings.filterwarnings('ignore')

```

```

datas = pd.read_csv("/kaggle/input/123345/complete_data_set.csv")
datas = datas.drop('Unnamed: 0', 1).values
X = datas[:, :-1]
X = min_max.fit_transform(X)
Y = datas[:, -1]
# 分层划分数据，采用五折家产验证
sfolder = StratifiedKFold(n_splits=5, shuffle=True, random_state=13)

# 激活函数
activation_list = ['relu']
# 求解器
solver_list = ['lbfgs']
alpha_list = [0.01 * i for i in range(0,1001,20)]

max_ = 0.0
string = None

for activation in activation_list:
    for solver in solver_list:
        for alpha in alpha_list:
            clf =
MLPClassifier(random_state=13,max_iter=200,activation=activation,solver=solver,alpha=alpha)
            scores = cross_val_score(clf, X, Y, cv=sfolder)
            avg_s = np.mean(np.array(scores))
            print("{} {} {} ->{}".format(activation, solver, alpha, avg_s))
            if max_ < avg_s:
                max_ = avg_s
                string = ("{} {} {} ->{}".format(activation, solver, alpha,
avg_s))

print("最好： " + string)

```

'''

- 运行最好效果参数
- ◦ relu lbfgs 4.6 ->0.3497148330358309

### 3、根据交叉验证得到的结果，所以选取最佳模型为随机森林

- 参数为：gini 8 log2

构建最佳模型代码：

'''key

```

# 随机森林
import numpy as np

```

```
import pandas as pd

# 随机森林
from sklearn.ensemble import RandomForestClassifier

# 获取数据
datas = pd.read_csv("complete_data_set.csv")
datas = datas.drop('Unnamed: 0',1).values
X = datas[:, :-1]
Y = datas[:, -1]

clf =
RandomForestClassifier(criterion='gini',max_depth=8,max_features='log2',n_estimators=120)

# 保存模型
from sklearn.externals import joblib

# 保存模型
joblib.dump(clf, 'RandomForest.pickle')
```

'''