

题目： 基于幸福度数据集的预测模型挖掘

一、人工智能对研究问题的定义和说明

二、结构框架技术和方法

- 1、总体思想
- 2、特征工程
- 3、模型选择
- 4、优化计算方法遗传算法
 - 4.1初始群体大小
 - 4.2适应度选择
 - 4.3交叉
 - 4.4变异

三、算法的代码实现（最后结果由于电脑性能差的问题没有跑出来，但是可以保证算法逻辑和代码实现没有问题）

- 1、读入数据、数据预处理
- 2、数据测试集和训练集划分及归一化
- 3、训练knn模型
- 4、训练BP神经网络模型
- 5、训练决策树模型
- 6、使用遗传算法，进行进化学习，得到每个决策分类模型中表现较好的7个参数模型，得到一共21个模型，最后进行集成投票的原则得到更好的分类结果。
 - 6.1、确定初始群体，进行初始化
 - 6.2、选择，使用准确率作为评分依据得到累计概率
 - 6.3、交叉，使用的是参数的交叉，随机根据参数的位置进行，交叉
 - 6.4、变异，通过对参数进行操作达到变异的目的
 - 6.5、进化大约200次最后得到21个较优的个体，最后得到综合的分类器
- 7、最后进行模型综合后，进行最后的评估

题目： 基于幸福度数据集的预测模型挖掘

一、人工智能对研究问题的定义和说明

幸福感是一个古老而深刻的话题，是人类世代追求的方向。与幸福感相关的因素成千上万、因人而异，大如国计民生，小如路边吃红薯，都会对幸福感产生影响。这些错综复杂的因素中，我们能找到其中的共性，一窥幸福感的要义吗？

通过中国综合社会调查（CGSS）2015年度调查问卷（居民问卷）得到了大量的数据集，该数据集包括了大量影响幸福度和幸福度因素的数据。

本文研究问题就是通过问卷调查数据，选取其中多组变量来得到一个人工智能预测模型来对人们的幸福度数值进行预测。

二、结构框架技术和方法

1、总体思想

首先，进行数据特征工程，然后选用的模型是多模型混合，即模型融合的方法来达到自己的预测模型，其中包括BP神经网络，k近邻方法，还有决策树模型。最后，进行模型参数的优化，使用进化计算中的遗传算法来进行参数优化。

2、特征工程

特征工程的主要工作是，探索数据集和处理数据集包括，对缺失值的处理，查看相关的文献或者资料，选择合适的属性特征，进行数据预处理等。

3、模型选择

模型选择为基本的多分类模型，包括BP神经网络，k近邻方法，还有决策树模型。达到一定的分类效果，进行模型融合。使用投票原则每个模型有一定的权重得到最后的分类结果。

4、优化计算方法遗传算法

使用遗传算法来优化，模型融合的权重参数以及各个模型的参数，最后得到自己的模型

4.1初始群体大小

选择合适的初始群体大小，合适的参数个体设计，作为初始基因，进行遗传。

4.2适应度选择

适应度即是，该个体更加适合遗传即分类效果越好，该个体适应度越大。因此使用分类准确率作为适应度函数，然后根据“轮盘赌”原则来获得新一轮的个体。

4.3交叉

个体与个体之间进行交叉，选择合适的位置，来进行参数之间的交叉，来产生一定数量的新个体，来达到获得更优秀基因的概率。

4.4变异

个体有一定的概率发生变异，产生新基因，寻求更优解。

三、算法的代码实现（最后结果由于电脑性能差的问题没有跑出来，但是可以保证算法逻辑和代码实现没有问题）

- 1、读入数据
- 2、先训练knn
- 3、训练bp神经网络
- 4、训练决策树
- 5、使用遗传算法,最后得到优化的参数
- 6、最后从最后的个体中,选取最优的21个个体,进行集成学习,进行幸福感预测
- 7、最后写报告整理资料、写课堂总结

1、读入数据、数据预处理

```
1 | # 导入相关的数据读入模块
2 | %matplotlib inline
3 | import pandas as pd
4 | import numpy as np
5 | import matplotlib.pyplot as plt
6 | import seaborn as sns
7 | import random
8 | plt.rcParams['font.sans-serif'] = ['KaiTi']
9 | plt.rcParams['font.serif'] = ['KaiTi']
```

```
1 | # 导入数据,解析survey_time列的值作为独立的日期列,指定字符集类型
2 | train = pd.read_csv('./data/happiness_train_complete.csv', parse_dates=['survey_time'], encoding='latin-1')
3 | train.head()
```

```
1 | .dataframe tbody tr th {
2 |     vertical-align: top;
3 | }
4 |
5 | .dataframe thead th {
6 |     text-align: right;
7 | }
```

| | id | happiness | survey_type | province | city | county | survey_time | gender | birth | nationality | ... | neighbor_f |
|---|----|-----------|-------------|----------|------|--------|---------------------|--------|-------|-------------|-----|------------|
| 0 | 1 | 4 | 1 | 12 | 32 | 59 | 2015-08-04 14:18:00 | 1 | 1959 | 1 | ... | 4 |
| 1 | 2 | 4 | 2 | 18 | 52 | 85 | 2015-07-21 15:04:00 | 1 | 1992 | 1 | ... | 3 |
| 2 | 3 | 4 | 2 | 29 | 83 | 126 | 2015-07-21 13:24:00 | 2 | 1967 | 1 | ... | 4 |
| 3 | 4 | 5 | 2 | 10 | 28 | 51 | 2015-07-25 17:33:00 | 2 | 1943 | 1 | ... | 3 |
| 4 | 5 | 4 | 1 | 7 | 18 | 36 | 2015-08-10 09:50:00 | 2 | 1994 | 1 | ... | 2 |

5 rows × 140 columns

```
1 | # 查看数据的形态
2 | train.shape
```

```
1 | (8000, 140)
```

```
1 | #查看所有的幸福度
2 | o=set(train['happiness'])
3 | print(o)
```

```
1 | {1, 2, 3, 4, 5, -8}
```

```
1 | # 处理异常值
2 | # 将happiness中的'-8'无法回答,改为'3'说不上幸福不幸福
3 | train['happiness'] = train['happiness'].replace(-8, 3)
```

```
1 | # 判断是否包含缺失值
2 | # isnull用于遍历dataframe每一个元素,将是空的元素置为True,将非空的元素置为False, sum用于计算每一列中为True元素的个数
3 | train.isnull().sum()
```

```

1 | id                0
2 | happiness         0
3 | survey_type       0
4 | province          0
5 | city              0
6 | county            0
7 | survey_time       0
8 | gender            0
9 | birth             0
10 | nationality        0
11 | religion           0
12 | religion_freq      0
13 | edu               0
14 | edu_other          7997
15 | edu_status         1120
16 | edu_yr             1972
17 | income             0
18 | political          0
19 | join_party         7176
20 | floor_area         0
21 | property_0         0
22 | property_1         0
23 | property_2         0
24 | property_3         0
25 | property_4         0
26 | property_5         0
27 | property_6         0
28 | property_7         0
29 | property_8         0
30 | property_other     7934
31 | ...
32 | m_political        0
33 | m_work_14          0
34 | status_peer        0
35 | status_3_before    0
36 | view               0
37 | inc_ability        0
38 | inc_exp            0
39 | trust_1            0
40 | trust_2            0
41 | trust_3            0
42 | trust_4            0
43 | trust_5            0
44 | trust_6            0
45 | trust_7            0
46 | trust_8            0
47 | trust_9            0
48 | trust_10           0
49 | trust_11           0
50 | trust_12           0
51 | trust_13           0
52 | neighbor_familiarity 0
53 | public_service_1    0
54 | public_service_2    0
55 | public_service_3    0
56 | public_service_4    0
57 | public_service_5    0
58 | public_service_6    0
59 | public_service_7    0
60 | public_service_8    0
61 | public_service_9    0
62 | Length: 140, dtype: int64

```

```

1 | # loc[] 方法用于实现将happiness=-8的行进行过滤，也就是保留happiness不为-8的行，因为幸福指数只有1-5，-8明显是一个错误的读入值
2 | train = train.loc[train['happiness'] != -8]

```

```

1 | train.shape

```

```

1 | (8000, 140)

```

```

1 | train['happiness'].value_counts()

```

```
1 4 4818
2 5 1410
3 3 1171
4 2 497
5 1 104
6 Name: happiness, dtype: int64
```

```
1 # 使用.dt.year将survey_time转换成year的时间
2 train['survey_time'] = train['survey_time'].dt.year
```

```
1 # 通过birth和survey时间计算出接受采访者的年龄
2 train['Age'] = train['survey_time']-train['birth']
```

```
1 del_list=['survey_time','birth']
2 print(train['Age'])
```

```
1 0 56
2 1 23
3 2 48
4 3 72
5 4 21
6 5 69
7 6 52
8 7 56
9 8 63
10 9 30
11 10 27
12 11 47
13 12 51
14 13 76
15 14 80
16 15 32
17 16 41
18 17 53
19 18 39
20 19 50
21 20 64
22 21 67
23 22 35
24 23 73
25 24 59
26 25 45
27 26 49
28 27 52
29 28 23
30 29 68
31 ..
32 7970 57
33 7971 20
34 7972 38
35 7973 39
36 7974 50
37 7975 43
38 7976 62
39 7977 67
40 7978 39
41 7979 31
42 7980 61
43 7981 69
44 7982 46
45 7983 62
46 7984 80
47 7985 45
48 7986 40
49 7987 52
50 7988 54
51 7989 84
52 7990 25
53 7991 43
54 7992 49
55 7993 61
56 7994 40
57 7995 34
58 7996 70
59 7997 48
60 7998 37
61 7999 24
62 Name: Age, Length: 8000, dtype: int64
```

```

1  '''
2  对dataframe中的数据进行简单的编码
3  小于16岁的为0；16-32的为1；32-48的为2；48-64的为3；64-80的为4；大于80的为5
4  '''
5  combine=[train]
6  for dataset in combine:
7      dataset.loc[dataset['Age']<=16, 'Age']=0
8      dataset.loc[(dataset['Age'] > 16) & (dataset['Age'] <= 32), 'Age'] = 1
9      dataset.loc[(dataset['Age'] > 32) & (dataset['Age'] <= 48), 'Age'] = 2
10     dataset.loc[(dataset['Age'] > 48) & (dataset['Age'] <= 64), 'Age'] = 3
11     dataset.loc[(dataset['Age'] > 64) & (dataset['Age'] <= 80), 'Age'] = 4
12     dataset.loc[ dataset['Age'] > 80, 'Age'] = 5

```

```

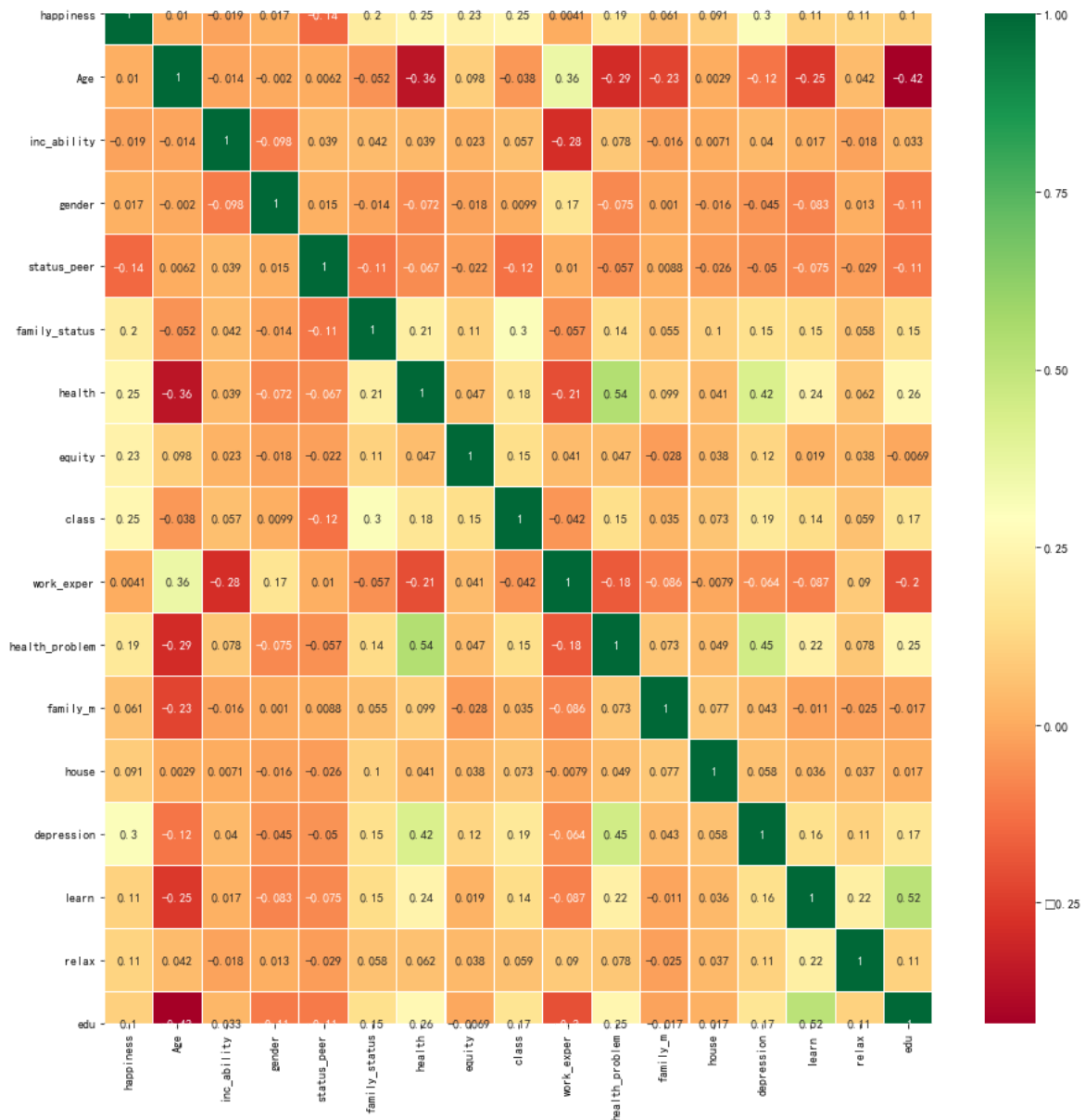
1  '''
2  绘制以下数据的热度图，
3  .corr()方法表示的是计算dataframe多个指标的相关系数矩阵，默认使用pearson计算方法
4  train[]表示传入热度图的数据，
5  annot（布尔类型），用于控制是否在个字中间标明数字，
6  cmap表示控制热度的渐变色，
7  linewidths表示每个单元格的线的宽度'''
8  sns.heatmap(train[['happiness', 'Age', 'inc_ability', 'gender', 'status_peer', 'family_status', 'health', 'equity', 'class', 'work_exper', 'health_problem', 'family_m', 'house', 'depression', 'learn', 'relax', 'edu']].corr(), annot=True, cmap='RdYlGn', linewidths=0.2) #data.corr()-->correlation matrix
9  fig=plt.gcf() #获取当前的图表和子图
10 fig.set_size_inches(15,15) #设置图像的密集度：设置图像的长和宽
11 plt.show()

```

```

1  C:\Users\Administrator\AppData\Roaming\Python\Python37\site-packages\matplotlib\backends\backend_agg.py:211: RuntimeWarning: Glyph 8722
missing from current font.
2      font.set_text(s, 0.0, flags=flags)
3  C:\Users\Administrator\AppData\Roaming\Python\Python37\site-packages\matplotlib\backends\backend_agg.py:180: RuntimeWarning: Glyph 8722
missing from current font.
4      font.set_text(s, 0, flags=flags)

```



```

1 '''
2 最后选择的特征为Age年龄,inc_ability收入是否合理,gender性别,status_peer与同龄人相比的收入情况,work_exper工作经历及情况,family_status家庭年收入情况,health身体健康状况,equity认为社会是否公平,class认为应该处于的社会阶层,health_problem影响健康的程度,family_m家庭人数,house拥有房产数量,depression压力沮丧程度,learn是否学习充电,relax休闲放松,edu教育程度
3 '''
4 features=
['Age','inc_ability','gender','status_peer','work_exper','family_status','health','equity','class','health_problem','family_m','house','depression','learn','relax','edu']

```

```

1 # 设置训练模型参数, target为模型输出的标签项, train为模型输入, test为测试集数据
2 # 特征工程完成得到最后用于训练和测试的数据集
3 train_target = train['happiness']
4 train = train[features]

```

2、数据测试集和训练集划分及归一化

```

1 # 进行数据随机划分模块
2 from sklearn.model_selection import train_test_split
3 # 进行数据标准化
4 from sklearn.preprocessing import StandardScaler
5 standardScaler = StandardScaler()
6 standardScaler.fit(train)
7 train_transform = standardScaler.transform(train)
8 X_train,X_test,Y_train,Y_test = train_test_split(train_transform,train_target,test_size=0.3)

```

3、训练knn模型

```

1 # 先引入进行评估的相关的包
2 from sklearn import metrics
3

```

```

1 from sklearn.neighbors import KNeighborsClassifier
2 test_knn = KNeighborsClassifier(n_neighbors = 5,p = 2,weights = 'distance')
3 test_knn.fit(X_train,Y_train)
4 Y_predict_train = test_knn.predict(X_test)
5 print(metrics.accuracy_score(Y_test,Y_predict_train,"%")

```

```

1 0.5591666666666667 %

```

4、训练BP神经网络模型

```

1 # MLPClassifier() BP神经网络
2 # 查看文章: https://blog.csdn.net/weixin\_42348202/article/details/100568469
3 from sklearn.neural_network import MLPClassifier

```

```

1 bp = MLPClassifier(hidden_layer_sizes=(50,50),max_iter=100)
2 bp.fit(X_train,Y_train)
3 Y_bp_predict = bp.predict(X_test)
4 print(metrics.accuracy_score(Y_test,Y_bp_predict,"%")

```

```

1 0.5870833333333333 %

```

```

1 E:\Anaconda3\lib\site-packages\sklearn\neural_network\multilayer_perceptron.py:566: ConvergenceWarning: Stochastic Optimizer: Maximum
iterations (100) reached and the optimization hasn't converged yet.
2 % self.max_iter, ConvergenceWarning)

```

5、训练决策树模型

```

1 # 导入需要的块
2 from sklearn.tree import DecisionTreeClassifier
3 tree = DecisionTreeClassifier(max_depth=7,min_samples_split=4)
4 tree.fit(X_train,Y_train)
5 Y_tree_predict = tree.predict(X_test)
6 print(metrics.accuracy_score(Y_tree_predict,Y_test),"%")

```

```

1 0.595 %

```

6、使用遗传算法，进行进化学习，得到每个决策分类模型中表现较好的7个参数模型，得到一共21个模型，最后进行集成投票的原则得到更好的分类结果。

```

1 # 定义一些全局变量，方便变异操作
2 # 所有模型对应的范围
3 # knn模型对应的范围
4 knn_n_neighbors_min = 5
5 knn_n_neighbors_max = 30
6 knn_p_min = 1
7 knn_p_max = 6
8 # bp神经网络模型参数对应范围
9 bp_iter_min = 100
10 bp_iter_max = 300
11 bp_x_min = 10
12 bp_x_max = 50
13 bp_y_min = 10
14 bp_y_max = 50
15 # 决策树模型参数对应的范围
16 tree_depth_min = 5
17 tree_depth_max = 15
18 tree_split_min = 2
19 tree_split_max = 8
20 #迭代的次数:
21 epoches = 200
22 # 相应的分类器集合和参数集合，设置为全局变量
23 knn_set = []
24 bp_set = []
25 tree_set = []
26 knn_param_set = []
27 bp_param_set = []
28 tree_param_set = []

```

6.1、确定初始群体，进行初始化

```

1 def init():
2     #knn初始化群体
3     for n_neighbors in range(5,31,2):
4         for p in range(1,6,1):
5             knn_set.append(KNeighborsClassifier(n_neighbors =n_neighbors ,p = p))
6             knn_param_set.append([n_neighbors,p])
7     for max_iter in range(100,300,20):
8         for x in range(10,51,5):
9             for y in range(10,51,5):

```

```

10         bp_set.append(MLPClassifier(hidden_layer_sizes=(x,y),max_iter=max_iter))
11         bp_param_set.append([max_iter,x,y])
12     for max_depth in range(5,15,1):
13         for min_samples_split in range(2,8,1):
14             tree_set.append(DecisionTreeClassifier(max_depth=max_depth,min_samples_split=min_samples_split))
15             tree_param_set.append([max_depth,min_samples_split])
16     return knn_set, bp_set, tree_set, knn_param_set, bp_param_set, tree_param_set

```

6.2、选择，使用准确率作为评分依据得到累计概率

```

1  # 得到评分
2  def classifier_score(X,Y,classifier):
3      classifier.fit(X, Y)
4      pred = classifier.predict(X)
5      return metrics.accuracy_score(y_true=Y, y_pred=pred)
6
7  # 计算每个个体的适应度,返回适应度圆盘，用于下一步轮盘赌原则用于选择个体
8  def adaption(X,Y,classifier_set):
9      score = []
10     for classifier in classifier_set:
11         score.append(classifier_score(X,Y,classifier))
12     sm = np.sum(score)
13     ada = score / sm
14     for i in range(1, len(ada)):
15         ada[i] = ada[i] + ada[i - 1]
16     return ada
17
18 # 根据轮盘进行选择,得到对应的分类器和对应的分类参数
19 def choose_classifiers(classifier_set,classifier_param_set,ada):
20     count = len(classifier_set)
21     result = []
22     result_param = []
23     for i in range(count):
24         r = random.random()
25         for j in range(len(ada)):
26             if r<= ada[j]:
27                 result.append(classifier_set[j])
28                 result_param.append(classifier_param_set[j])
29                 break;
30     return result,result_param

```

6.3、交叉，使用的是参数的交叉，随机根据参数的位置进行，交叉

```

1  # 两个分类器参数之间的交叉，分类型进行不用类别的交叉，
2  def cross(index1,index2,type,param_set):
3      if (type == "knn") or (type == "tree"):
4          #随机产生0-1的数进行交叉
5          index_temp = random.randint(0,1)
6          data_temp = param_set[index1][index_temp]
7          param_set[index1][index_temp] = param_set[index2][index_temp]
8          param_set[index2][index_temp] = data_temp
9      else:
10         # 随机产生0-2的树进行交叉
11         index_temp = random.randint(0,2)
12         data_temp = param_set[index1][index_temp]
13         param_set[index1][index_temp] = param_set[index2][index_temp]
14         param_set[index2][index_temp] = data_temp
15     return
16
17 # 整个集合的交叉也要传进来类型是什么，对于该集合的参数，每间隔三个个体进行随机交叉
18 def all_cross(classifier_set,type):
19     set_len = len(classifier_set)
20     # 最后交叉完的集合
21     for i in range(1,set_len,3):
22         cross(i,i-1,type,classifier_set)
23     return

```

6.4、变异，通过对参数进行操作达到变异的目的

```

1  # 分类型进行不同类别的变异，变异的范围还要进行自己的确认不能超出正常范围，变异率默认设置为0.2
2  # 就是简单的+1或者-1操作，当达到最大值时只能减，最小值时只能加，其他随意加减
3  def variation(classifier_set,type,probability):
4      count = int(len(classifier_set)*probability)#变异的个数
5      while count > 0:
6          index_tep = random.randint(0,len(classifier_set))#每个分类器变异的位置
7          operation = random.randint(0,1)#0减1加
8          if type == 'knn':
9              index_param = random.randint(0,1)#参数变异的位置
10             if operation == 0:
11                 if index_param == 0:
12                     if classifier_set[index_tep][0] == knn_n_neighbors_min:
13                         classifier_set[index_tep][0] = classifier_set[index_tep][0] +1
14                     else:
15                         classifier_set[index_tep][0] = classifier_set[index_tep][0] -1
16             else:

```



```

17         if classifier_set[index_tep][1] == knn_p_min:
18             classifier_set[index_tep][1] = classifier_set[index_tep][1] + 1
19         else:
20             classifier_set[index_tep][1] = classifier_set[index_tep][1] - 1
21     else: #运算是加的情况
22         if index_param == 0:
23             if classifier_set[index_tep][0] == knn_n_neighbors_max:
24                 classifier_set[index_tep][0] = classifier_set[index_tep][0] - 1
25             else:
26                 classifier_set[index_tep][0] = classifier_set[index_tep][0] + 1
27         else:
28             if classifier_set[index_tep][1] == knn_p_max:
29                 classifier_set[index_tep][1] = classifier_set[index_tep][1] - 1
30             else:
31                 classifier_set[index_tep][1] = classifier_set[index_tep][1] + 1
32 elif type == "tree":
33     index_param = random.randint(0,1) #参数变异的位置
34     if operation == 0:
35         if index_param == 0:
36             if classifier_set[index_tep][0] == tree_depth_min:
37                 classifier_set[index_tep][0] = classifier_set[index_tep][0] + 1
38             else:
39                 classifier_set[index_tep][0] = classifier_set[index_tep][0] - 1
40         else:
41             if classifier_set[index_tep][1] == tree_split_min:
42                 classifier_set[index_tep][1] = classifier_set[index_tep][1] + 1
43             else:
44                 classifier_set[index_tep][1] = classifier_set[index_tep][1] - 1
45     else: #运算是加的情况
46         if index_param == 0:
47             if classifier_set[index_tep][0] == tree_depth_max:
48                 classifier_set[index_tep][0] = classifier_set[index_tep][0] - 1
49             else:
50                 classifier_set[index_tep][0] = classifier_set[index_tep][0] + 1
51         else:
52             if classifier_set[index_tep][1] == tree_split_max:
53                 classifier_set[index_tep][1] = classifier_set[index_tep][1] - 1
54             else:
55                 classifier_set[index_tep][1] = classifier_set[index_tep][1] + 1
56 else: #是神经网络模型的时候
57     index_param = random.randint(0,2) #参数变异的位置
58     if operation == 0:
59         if index_param == 0:
60             if classifier_set[index_tep][0] == bp_iter_min:
61                 classifier_set[index_tep][0] = classifier_set[index_tep][0] + 1
62             else:
63                 classifier_set[index_tep][0] = classifier_set[index_tep][0] - 1
64         elif index_param == 1:
65             if classifier_set[index_tep][1] == bp_x_min:
66                 classifier_set[index_tep][1] = classifier_set[index_tep][1] + 1
67             else:
68                 classifier_set[index_tep][1] = classifier_set[index_tep][1] - 1
69         else:
70             if classifier_set[index_tep][2] == bp_y_min:
71                 classifier_set[index_tep][2] = classifier_set[index_tep][2] + 1
72             else:
73                 classifier_set[index_tep][2] = classifier_set[index_tep][2] - 1
74     else: #运算是加的情况
75         if index_param == 0:
76             if classifier_set[index_tep][0] == bp_iter_max:
77                 classifier_set[index_tep][0] = classifier_set[index_tep][0] - 1
78             else:
79                 classifier_set[index_tep][0] = classifier_set[index_tep][0] + 1
80         elif index_param == 1:
81             if classifier_set[index_tep][1] == bp_x_max:
82                 classifier_set[index_tep][1] = classifier_set[index_tep][1] - 1
83             else:
84                 classifier_set[index_tep][1] = classifier_set[index_tep][1] + 1
85         else:
86             if classifier_set[index_tep][2] == bp_y_max:
87                 classifier_set[index_tep][2] = classifier_set[index_tep][2] - 1
88             else:
89                 classifier_set[index_tep][2] = classifier_set[index_tep][2] + 1
90     count = count - 1

```

```

1  # 更新相应的分类器参数后, 继续根据参数, 得到相应的分类器集合
2  def to_get_classifiers_set():
3      new_knn_set = []
4      new_bp_set = []
5      new_tree_set = []
6      for temp in knn_param_set:
7          new_knn_set.append(KNeighbor, sClassifier(n_neighbors = temp[0], p = temp[1]))
8      for temp1 in bp_param_set:
9          new_bp_set.append(MLPClassifier(hidden_layer_sizes=(temp1[1], temp1[2]), max_iter=temp[0]))
10     for temp2 in tree_param_set:
11         new_tree_set.append(DecisionTreeClassifier(max_depth=temp2[0], min_samples_split=temp2[1]))
12     return new_knn_set, new_bp_set, new_tree_set

```

6.5、进化大约200次最后得到21个较优的个体, 最后得到综合的分类器

```

1  # 迭代200次得到每个类别进化的个体集合, 从每个集合选择准确率最高的7个
2  def to_get_best(knn_set, bp_set, tree_set):
3      knn_best_set = set()
4      bp_best_set = set()
5      tree_best_set = set()
6      min_score = 1 # 记录当前集合最小分数
7      min_classifier = 0 # 记录最小评估分数的分类器
8      for i in range(0, 7):
9          for temp_classifier in knn_set:
10             if len(knn_best_set) <= 7:
11                 temp_score = classifier_score(X_train, X_test, temp_classifier)
12                 if temp_score < min_score:
13                     min_classifier = temp_classifier
14                     knn_best_set.add(temp_classifier)
15             else:
16                 temp_score = classifier_score(X_train, X_test, temp_classifier)
17                 if temp_score < min_score:
18                     knn_best_set.remove(min_classifier)
19                     knn_best_set.add(temp_classifier)
20                 min_score = temp_score
21         for temp_classifier in bp_set:
22             if len(bp_best_set) <= 7:
23                 temp_score = classifier_score(X_train, X_test, temp_classifier)
24                 if temp_score < min_score:
25                     min_classifier = temp_classifier
26                     bp_best_set.add(temp_classifier)
27             else:
28                 temp_score = classifier_score(X_train, X_test, temp_classifier)
29                 if temp_score < min_score:
30                     bp_best_set.remove(min_classifier)
31                     bp_best_set.add(temp_classifier)
32                 min_score = temp_score
33         for temp_classifier in tree_set:
34             if len(knn_best_set) <= 7:
35                 temp_score = classifier_score(X_train, X_test, temp_classifier)
36                 if temp_score < min_score:
37                     min_classifier = temp_classifier
38                     tree_best_set.add(temp_classifier)
39             else:
40                 temp_score = classifier_score(X_train, X_test, temp_classifier)
41                 if temp_score < min_score:
42                     tree_best_set.remove(min_classifier)
43                     tree_best_set.add(temp_classifier)
44                 min_score = temp_score
45         the_end_set = set.union(knn_best_set, bp_best_set, tree_best_set)
46     return the_end_set
47

```

```

1  ##### 6、根据现在所有的模型, 进行投票集成学习得到最后所有元素的分类, 并且返回准确率
2  def to_ensemble(end_classifiers_set):
3      y_label_predict = []
4      for temp_data in X_test:
5          predict_set = {}
6          for classifiers in end_classifiers_set:
7              classifiers.fit(X_train, Y_train)
8              predict_set.append(classifiers.predict(temp_data))
9              if classifiers not in predict_set.keys():
10                 predict_set[classifiers] = 1
11             else:
12                 predict_set[classifiers] = predict_set[classifiers] + 1
13             label = max(predict_set, key=predict_set.get)
14             y_label_predict.append(label)
15     return metrics.accuracy_score(y_label_predict, Y_test)

```

7、最后进行模型综合后, 进行最后的评估

```

1  # 利用最后得到的最优的21个个体最终得到最后的模型
2  # 主函数入口
3  def main():

```

```
4      #初始化
5      global epochs
6      knn_set,bp_set,tree_set,knn_param_set,bp_param_set,tree_param_set = init()
7      while epochs > 0 :
8          #选择
9          ada1 = adaption(X_train,Y_train,knn_set)
10         knn_set,knn_param_set = choose_classifiers(knn_set,knn_param_set,ada1)
11         ada2 = adaption(X_train,Y_train,bp_set)
12         bp_set,bp_param_set = choose_classifiers(bp_set,bp_param_set,ada2)
13         ada3 = adaption(X_train,Y_train,tree_set)
14         tree_set,tree_param_set = choose_classifiers(tree_set,tree_param_set,ada3)
15         #交叉
16         all_cross(knn_param_set,"knn")
17         all_cross(bp_param_set,"bp")
18         all_cross(tree_param_set,"tree")
19         #变异
20         variation(knn_param_set,"knn",0.2)
21         variation(bp_param_set,"bp",0.2)
22         variation(tree_param_set,"tree",0.2)
23         #更新分类集合
24         knn_set,bp_set,tree_set = to_get_classifiers_set()
25         epochs = epochs -1
26     #得到综合分类器
27     the_end_set = to_get_best(knn_set,bp_set,tree_set)
28     accuracy = to_ensemble(the_end_set)
29     print("准确率是: ",accuracy)
```

1 | main()