

# AdaBoost代码解析

准确率基本为100%

```
1  """
2  * @author  孟子喻
3  * @time    2021.5.16
4  * @file    classify_iris.py
5  *          adaboost.py
6  """
7  import numpy as np
8  import pandas as pd
9  from sklearn.datasets import load_iris
10 from adaboost import AdaBoost
11 import matplotlib.pyplot as plt
12 import csv
13 from sklearn.model_selection import train_test_split
14
15 # 鸢尾花(iris)数据集
16 # 数据集内包含 3 类共 150 条记录，每类各 50 个数据，
17 # 每条记录都有 4 项特征：花萼长度、花萼宽度、花瓣长度、花瓣宽度，
18 # 可以通过这4个特征预测鸢尾花卉属于（iris-setosa, iris-versicolour, iris-
19 #  virginica）中的哪一品种。
20 # 这里只取前100条记录，两项特征，两个类别。
21
22 def create_data():
23     iris = load_iris()
24     df = pd.DataFrame(iris.data, columns=iris.feature_names)
25     df['label'] = iris.target
26     df.columns = ['sepal length', 'sepal width', 'petal length', 'petal
27 width', 'label']
28     data = np.array(df.iloc[:100, [0, 1, -1]])
29     for i in range(len(data)):
30         if data[i, -1] == 0:
31             data[i, -1] = -1
32     # print(data)
33     return data
34
35 def calc_acc(y, y_hat):
36     idx = np.where(y_hat == 1)
37     TP = np.sum(y_hat[idx] == y[idx])
38     idx = np.where(y_hat == -1)
39     TN = np.sum(y_hat[idx] == y[idx])
40     return float(TP + TN)/len(y)
41
42 if __name__ == "__main__":
43     trainset = create_data()
44
45     features, labels = trainset[:, :2], trainset[:, -1]
46     x = features
47     y = labels
48
```

```

49     ada = AdaBoost()
50     clf_num = ada.fit(x, y) # 返回使用的弱分类器数量
51     print(clf_num)
52     ada.show_para() # 显示一些参数
53
54     # 进行预测
55     y_hat = ada.output_final_pred(x, clf_num)
56
57     acc = calc_acc(y, y_hat) # 计算准确率
58     print("Accuracy: ", str(100* acc), "%")
59
60     plt.figure()
61
62     for i in range(clf_num):
63         if ada.classifier[i].decision_feature == 0:
64             plt.plot([ada.classifier[i].decision_threshold,
65                        ada.classifier[i].decision_threshold], [0, 10])
66         else:
67             plt.plot([0, 10], [ada.classifier[i].decision_threshold,
68                                ada.classifier[i].decision_threshold])
69
70     for item in trainset:
71         if item[-1] == 1:
72             plt.scatter(item[0], item[1], s=30, marker='+', color='blue')
73         else:
74             plt.scatter(item[0], item[1], s=30, marker='_', color='red')
75     plt.xlim(3, 8)
76     plt.ylim(1, 6)
77     plt.show()

```

这一段代码和我之前的代码一样，用于把数据集分类，然后输入模型

下面是AdaBoost的代码详解

初始化模型

```

1 class AdaBoost:
2     def __init__(self):
3         self.w = {} # 权重weight
4         self.classifier = {} # 分类器字典，存储用到的多个弱分类器
5         self.alpha = {} # 每个分类器的影响权重
6         self.pred = {} # 每个分类器预测的值
7         self.classifier_num = 0 # 弱分类器数量

```

以下是具体的训练代码，每一句的作用基本都放在注释里了

```

1 # 训练
2 def fit(self, x, y, max_clf_num=15):
3
4     pred = {}
5     data_size = len(y) # 直接将所有实例一个batch进行训练，所以要经常用到一个
6     batch的大小
7     for iter in range(max_clf_num):

```

```

8         # 在允许的最大分类器数量范围内，对每个弱分类器进行初始化
9         self.w.setdefault(iter)          # 每一个变量的权重，不同弱分类器
对应的权重不一样
10        self.classifier.setdefault(iter)   # 分类器字典
11        self.alpha.setdefault(iter)        # 该分类器的alpha值
12        self.pred.setdefault(iter)         # 该分类器的预测值
13        # setdefault为检查iter是否存在于字典的keys中，若不存在则新建key并将
value设置为None
14
15        # 进行迭代
16        for iter in range(max_clf_num):
17            # 第一次迭代时初始化实例权重w，此时每个实例对分类器影响相同，所以设置为
1/data_size
18            if iter == 0:
19                self.w[iter] = np.ones(data_size) / data_size
20                self.w[iter] = self.w[iter].reshape([data_size, 1])
21            # 如果不是第一次迭代，则更新weight
22            else:
23                self.w[iter] = self.cal_w(iter, y, pred[iter - 1])
24            # pred[iter-1]是计算出的新的预测值
25
26            # 使用弱分类器进行分类
27            self.classifier[iter] = LineClassifier()
28            self.classifier[iter].fit(x, y, self.w[iter])
29            pred[iter] = self.classifier[iter].pred      # 使用弱分类器进行一轮预
测
30
31            # 计算误差
32            error = self.cal_e(y, pred[iter], self.w[iter])
33            # 计算alpha
34            self.alpha[iter] = self.cal_alpha(error)
35            # 计算最终预测值（每次更新）
36            final_predict = self.cal_final_pred(iter, self.classifier,
data_size)
37
38            print('iteration:%d' % (iter + 1))
39            self.classifier_num = iter
40
41            # 输出最终结果
42            if self.cal_final_e(y, final_predict) == 0 or error == 0:
43                print('cal_final_predict:%s' % (final_predict))
44                break
45                print('self.decision_key=%s' %
(self.classifier[iter].decision_key))
46                print('self.decision_feature=%d' %
(self.classifier[iter].decision_feature))
47                print('decision_threshold=%f' %
(self.classifier[iter].decision_threshold))
48                print('error:%f alpha:%f' % (error, self.alpha[iter]))
49                print('\n')
50            # 输出弱分类器数量
51            return self.classifier_num

```

主函数中用到的一些其他函数，基本都对应课本公式

```

1    def cal_e(self, y, pred, w):
2        ret = 0

```

```

3         for i in range(len(y)):
4             if y[i] != pred[i]:
5                 ret += w[i]
6         return ret
7
8     # 迭代更新alpha
9     def cal_alpha(self, e):
10         if e == 0:
11             return 10000
12         elif e == 0.5:
13             return 0.001
14         else:
15             return 0.5 * np.log((1 - e) / e) # 对应课本公式8.2，返回Gm的系数
16
17     # 计算最终的预测值
18     def cal_final_pred(self, i, classifier, data_size):
19         ret = np.array([0.0] * data_size)
20         for j in range(i + 1):
21             ret += self.alpha[j] * classifier[j].pred
22         return np.sign(ret)

```

## 以下是弱分类器

### 初始化部分

```

1 class LineClassifier:
2     def __init__(self):
3         self.w = None # 传入的每个实例的权重
4         self.decision_key = None # 取大于阈值的一侧还是小于的一侧（因为是划定
    正方形一个区域，所以要有方向）
5         self.decision_feature = None # 哪个特征起到分类作用x或y
6         self.decision_threshold = None # 分类的阈值
7         self.pred = None # 该分类器本次的预测值

```

### 训练部分

```

1 # 训练
2 def fit(self, x, y, w):
3     self.w = w # w是特征权重
4     dic = self.cal_dic(x) # 计算用每个实例的每个特征做阈值时，所有实例点在阈值两
    侧的分类情况
5     e_dic = self.cal_e_dic(y, dic) # 计算生成对应的error字典
6     e_min, self.decision_key, self.decision_feature, e_min_i =
    self.cal_e_min(e_dic) # 计算最小的error
7     self.decision_threshold = x[e_min_i, self.decision_feature] # 把最
    小error对应的阈值取出来输出
8     self.pred = dic[self.decision_key][self.decision_feature][e_min_i]
    # 用本次训练出的弱分类器进行一次预测

```

### 统计数据以便于后续阈值选取部分

```

1 def cal_dic(self, x):
2     ret_gt = {} # gt代表greater than
3     for i in range(x.shape[1]):
4         # 此处为取x或y,记为特征i
5         ret_gt[i] = []

```

```

6         for j in range(X.shape[0]):
7             # 此处为取一个实例，记为实例j
8             temp_threshold = X[j, i]    # 将取到的某个实例j的特征i作为阈值
9             temp_line = []
10            for k in range(X.shape[0]):    # 继续取所有实例j的特征i，如果超
过了阈值，则标记为1存入temp_line
11                if X[k, i] >= temp_threshold:
12                    temp_line.append(1)
13                    # 如果
14            else:                                # 如果没超过阈值，则标记为-1存入
temp_line
15                temp_line.append(-1)
16            ret_gt[i].append(temp_line)    # 将取到的与阈值的比较情况存入
ret_gt[i]
17
18            ret_lt = {}    # lt代表lower than
19            for i in range(X.shape[1]):
20                ret_lt[i] = []
21                for j in range(X.shape[0]):
22                    temp_threshold = X[j, i]
23                    temp_line = []
24                    for k in range(X.shape[0]):
25                        if X[k, i] <= temp_threshold:
26                            temp_line.append(1)
27                        else:
28                            temp_line.append(-1)
29                    ret_lt[i].append(temp_line)
30            ret = {'gt': ret_gt, 'lt': ret_lt}    # 是一个字典，里面存储了每一个实例的每
一个特征作为阈值时，有几个点大于或小于它
31            return ret

```