

机器学习学习报告

周珊琳

上海电力大学

上海, 2019-07-23, 中国

第一节 算法概述

算法名称: AdaBoost

算法原理: Boosting, 也称为增强学习或提升法, 是一种重要的集成学习技术, 能够将预测精度仅比随机猜度略高的弱学习器增强为预测精度高的强学习器, 这在直接构造强学习器非常困难的情况下, 为学习算法的设计提供了一种有效的新思路和新方法。作为一种元算法框架, Boosting 几乎可以应用于所有目前流行的机器学习算法以进一步加强原算法的预测精度, 应用十分广泛, 产生了极大的影响。而 AdaBoost 正是其中最成功的代表, 被评为数据挖掘十大算法之一。在 AdaBoost 提出至今的十几年间, 机器学习领域的诸多知名学者不断投入到算法相关理论的研究中去, 扎实的理论为 AdaBoost 算法的成功应用打下了坚实的基础。

第二节 算法设计

2.1 算法流程

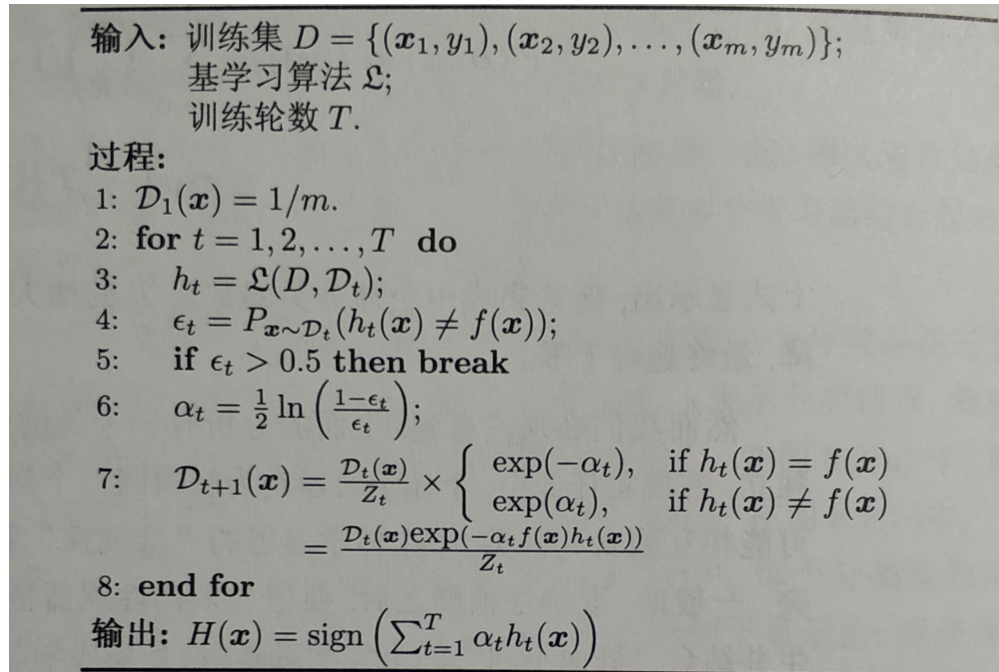


图 1: AdaBoost 算法流程 (引自西瓜书 174 页图 8.3)

2.2 核心代码

源代码 1:

```
1  # -*- coding: utf-8 -*-
2  """
3  Created on Thu Jul 23 15:55:05 2019
4
5  @author: zsl
6  """
7
8  import numpy as np
9  from sklearn.datasets import load_iris, load_wine
10 from sklearn import preprocessing
11 import matplotlib.pyplot as plt
12
13 def process_data(data):
14     min_max_scaler = preprocessing.MinMaxScaler()
15     return min_max_scaler.fit_transform(data)
16
17 #iris, wine
18 def loadData(data, y1):
19     #load
20     dataSet = data.data
21     target = data.target
22     #shuffle
23     num_example = dataSet.shape[0]
24     array = np.arange(num_example)
25     np.random.shuffle(array)
26     dataSet = dataSet[array]
27     target = target[array]
28     #transform
29     classLabels = list(target)
30     dataArr = process_data(dataSet)
31     for i in range(len(classLabels)):
32         if classLabels[i] == y1:
33             classLabels[i] = 1.0
34         else:
35             classLabels[i] = -1.0
36     return dataArr, classLabels
37
38 #car
39 def loadDataSet(fileName):
40     dataMat = []; classLabels = []
41     fr = open(fileName)
42     for line in fr.readlines():
43         lineArr = line.strip().split(',')
44         # print(lineArr)
45         if lineArr[0] == 'vhigh':
46             lineArr[0] = 1
47         if lineArr[0] == 'high':
48             lineArr[0] = 2
49         if lineArr[0] == 'med':
50             lineArr[0] = 3
```

```

51     if lineArr[0] == 'low':
52         lineArr[0] = 4
53     if lineArr[1] == 'vhigh':
54         lineArr[1] = 1
55     if lineArr[1] == 'high':
56         lineArr[1] = 2
57     if lineArr[1] == 'med':
58         lineArr[1] = 3
59     if lineArr[1] == 'low':
60         lineArr[1] = 4
61     if lineArr[2] == '2':
62         lineArr[2] = 1
63     if lineArr[2] == '3':
64         lineArr[2] = 2
65     if lineArr[2] == '4':
66         lineArr[2] = 3
67     if lineArr[2] == '5more':
68         lineArr[2] = 4
69     if lineArr[3] == '2':
70         lineArr[3] = 1
71     if lineArr[3] == '4':
72         lineArr[3] = 2
73     if lineArr[3] == 'more':
74         lineArr[3] = 3
75     if lineArr[4] == 'small':
76         lineArr[4] = 1
77     if lineArr[4] == 'med':
78         lineArr[4] = 2
79     if lineArr[4] == 'big':
80         lineArr[4] = 3
81     if lineArr[5] == 'low':
82         lineArr[5] = 1
83     if lineArr[5] == 'med':
84         lineArr[5] = 2
85     if lineArr[5] == 'high':
86         lineArr[5] = 3
87     dataMat.append([float(lineArr[0]), float(lineArr[1]), float(lineArr[2]),
88                     float(lineArr[3]), float(lineArr[4]),
89                     float(lineArr[5])])
90
91
92     if lineArr[6] == 'unacc':
93         lineArr[6] = 1.0
94     # elif lineArr[6] == 'acc':
95     #     lineArr[6] = 1
96     # elif lineArr[6] == 'good':
97     #     lineArr[6] = 2
98     else:
99         lineArr[6] = -1.0
100
101     classLabels.append(float(lineArr[6]))
102
103     return np.array(dataMat, dtype="float64"), classLabels
104

```

```

105 def stumpClassify(dataMatrix,col,threshVal,threshIneq):
106     retArray =np.ones((np.shape(dataMatrix)[0],1))
107     if threshIneq == 'lt':
108         retArray[dataMatrix[:,col] <=threshVal] =-1.0
109     else:
110         retArray[dataMatrix[:,col] >threshVal] =-1.0
111     return retArray
112
113 def buildStump(dataArr,classLabels,W):
114     dataMatrix =np.mat(dataArr)
115     classMatrix =np.mat(classLabels).T
116     m,n =np.shape(dataMatrix)
117     numSteps =10.0
118     bestStump={}
119     bestClasEst =np.mat(np.zeros((m,1)))
120     minError =np.inf
121     for i in range(n):
122         colMax =dataMatrix[:,i].max()
123         colMin =dataMatrix[:,i].min()
124         stepSize =(colMax-colMin)/numSteps
125         for j in range(-1,int(numSteps)+1):
126             for inequal in ['lt','gt']:
127                 threshVal =(colMin +float(j)*stepSize)
128                 predictedVals =stumpClassify(dataMatrix,i,threshVal,inequal)
129                 # print('predictedVals:=====','predictedVals)
130                 errArr =np.mat(np.ones((m,1)))
131                 # print('errArr',errArr)
132                 errArr[predictedVals ==classMatrix] =0
133                 weightedError =W.T*errArr
134                 # print("split: col %d, thresh %.2f, thresh inequal: %s, the weighted error is %.3f" % (i,
135                                                         #                                     threshVal, inequal, weightedError))
136
137                 if weightedError <minError:
138                     minError =weightedError
139                     bestClasEst =predictedVals.copy()
140                     bestStump['col']=i
141                     bestStump['thresh']=threshVal
142                     bestStump['ineq'] =inequal
143             return bestStump,minError,bestClasEst
144
145 def adaBoostTrainDS(dataArr,classLables,T=30):
146     weakClassArr =[]
147     errorList =[]
148     m =np.shape(dataArr)[0]
149     W =np.mat(np.ones((m,1))/m)
150     aggClassEst =np.mat(np.zeros((m,1)))
151     for i in range(T):
152         bestStump,error,classEst=buildStump(dataArr,classLables,W)
153         alpha =float(0.5*np.log((1.0-error)/max(error,1e-16)))
154         bestStump['alpha'] =alpha
155         weakClassArr.append(bestStump)
156         expon =np.multiply(-1*alpha*np.mat(classLables).T,classEst)
157         W =np.multiply(W,np.exp(expon))
158         W =W/W.sum()
159         aggClassEst +=alpha*classEst

```

```

158 #     print(aggClassEst)
159     aggErrors = np.multiply(np.sign(aggClassEst) != np.mat(classLabels).T, np.ones((m, 1)))
160 #     print('=====', aggErrors)
161     errorRate = aggErrors.sum() / m
162 #     print('total error', errorRate)
163     errorList.append(errorRate)
164     if errorRate == 0.0:
165         break
166     return weakClassArr, aggClassEst, errorList
167
168 def adaClassify(datToClass, classifierArr):
169     dataMatrix = np.mat(datToClass)
170     m = np.shape(dataMatrix)[0]
171     aggClassEst = np.mat(np.zeros((m, 1)))
172     for i in range(len(classifierArr)):
173         classEst = stumpClassify(dataMatrix, classifierArr[i]['col'],
174                                   classifierArr[i]['thresh'],
175                                   classifierArr[i]['ineq'])
176         aggClassEst += classifierArr[i]['alpha'] * classEst
177 #     print(np.sign(aggClassEst))
178     return np.sign(aggClassEst)
179
180 def classifytest(testDataSet, classifierArr, target_test):
181     """
182     计算准确率
183     """
184     i = 0
185     cnt = 0
186     for testVec in testDataSet:
187         pre = adaClassify(testDataSet, classifierArr)
188         if (pre == target_test[i]).all():
189             cnt += 1
190         i += 1
191 #     print('cnt', cnt)
192     return cnt / len(target_test)
193
194 ##iris
195 data = load_iris()
196 dataArr1, classLabels1 = loadDataSet(data, 0)
197
198 #wine
199 data = load_wine()
200 dataArr2, classLabels2 = loadDataSet(data, 3)
201
202 ##car
203 dataArr3, classLabels3 = loadDataSet('car.data')
204 #W = np.mat(np.ones((np.shape(dataArr)[0], 1)) / np.shape(dataArr)[0])
205 #bestStump, minError, bestClassEst = buildStump(dataArr, classLabels, W)
206 #classifierArr, aggClassEst, errorList = adaBoostTrainDS(dataArr1, classLabels1, 50)
207
208 num_example = dataArr2.shape[0]
209 sample = np.int(num_example * 0.9)
210 x_train = dataArr2[:sample]
211 y_train = classLabels2[:sample]

```

```

212 x_test =dataArr2[sample:]
213 y_test =classLabels2[sample:]
214 res=[]
215 for j in range(10):
216     pri =[]
217     tmp =0
218     prob =0.1
219     for i in range(10):
220         classifierArr,aggClassEst,errorList =adaBoostTrainDS(x_train[tmp:np.int(sample*prob)],y_train[
                tmp:np.int(sample*prob)],50)
221         pre=classifytest(x_test, classifierArr,y_test)
222         pri.append(pre)
223         tmp=np.int(sample*prob)+1
224         prob +=0.1
225     res.append(sum(pri)/len(pri))
226     print(sum(pri)/len(pri))
227 print('wine AdaBoost, 10次十折交叉验证结果:',(sum(res)/len(res)))
228
229
230
231 #plt.rcParams['font.sans-serif']=['SimHei'] #用来正常显示中文标签
232 #plt.rcParams['axes.unicode_minus']=False #用来正常显示负号
233 #plt.figure()
234 #ln1, = plt.plot(errorList1,linestyle='dashed',linewidth=0.5,color='red',marker='.',)
235 #ln2, = plt.plot(errorList2,linestyle='dashed',linewidth=0.5,color='b',marker='.',)
236 #ln3, = plt.plot(errorList3,linestyle='dashed',linewidth=0.5,color='g',marker='.',)
237 #ln4, = plt.plot(errorList4,linestyle='dashed',linewidth=0.5,color='yellow',marker='.',)
238 #plt.ylim(0,0.3)
239 #plt.legend(handles=[ln1,ln2,ln3,ln4], labels=['unacc', 'acc','good','vgood'],
240 #         loc='upper right')
241 #plt.ylabel('errorRate')
242 #plt.title('训练50次过程中数据集car错误率的变化情况')
243 #plt.show()

```

第三节 选用数据

iris 行数: 150 列数: 5

列属性及取值:

- 1) 萼片长度 cm, 数值型
- 2) 萼片宽度 cm, 数值型
- 3) 花瓣长度 cm, 数值型
- 4) 花瓣宽度 cm 数值型

类别:

Iris Setosa

Iris Versicolour

Iris Virginica

wine, 行数: 178, 列数: 13

属性:

- 1) Alcohol
- 2) Malic acid
- 3) Ash
- 4) Alcalinity of ash
- 5) Magnesium
- 6) Total phenols
- 7) Flavanoids
- 8) Nonflavanoid phenols
- 9) Proanthocyanins
- 10) Color intensity
- 11) Hue
- 12) OD280/OD315 of diluted wines
- 13) Proline

类别:

Alcohol 1, 2, 3

car, 行数: 1728, 列数: 6

列属性及取值:

- 1) buying: vhigh, high, med, low.
- 2) maint: vhigh, high, med, low.
- 3) doors: 2, 3, 4, 5more.
- 4) persons: 2, 4, more.
- 5) lugboot: small, med, big.
- 6) safety: low, med, high.

类别:

unacc, acc, good, vgood

第四节 实验结果展示

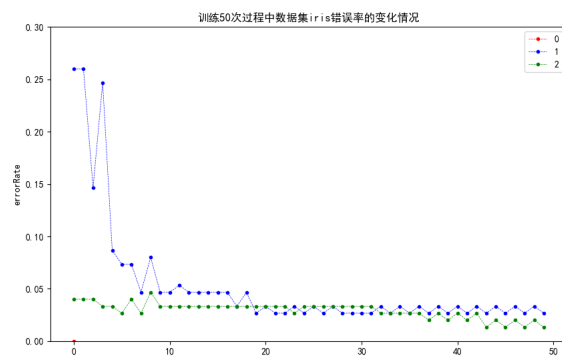


图 2: iris 类标号选择对比图

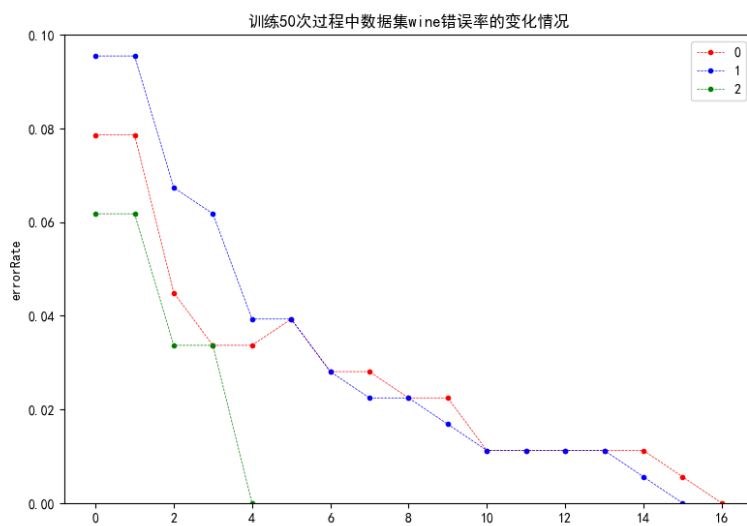


图 3: wine 类标号选择对比图

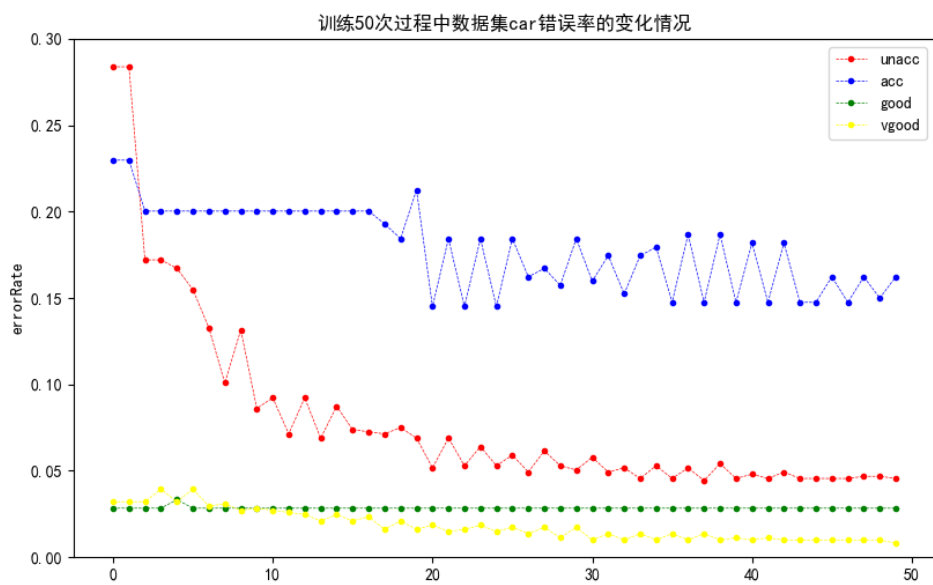


图 4: car 类标号选择对比图

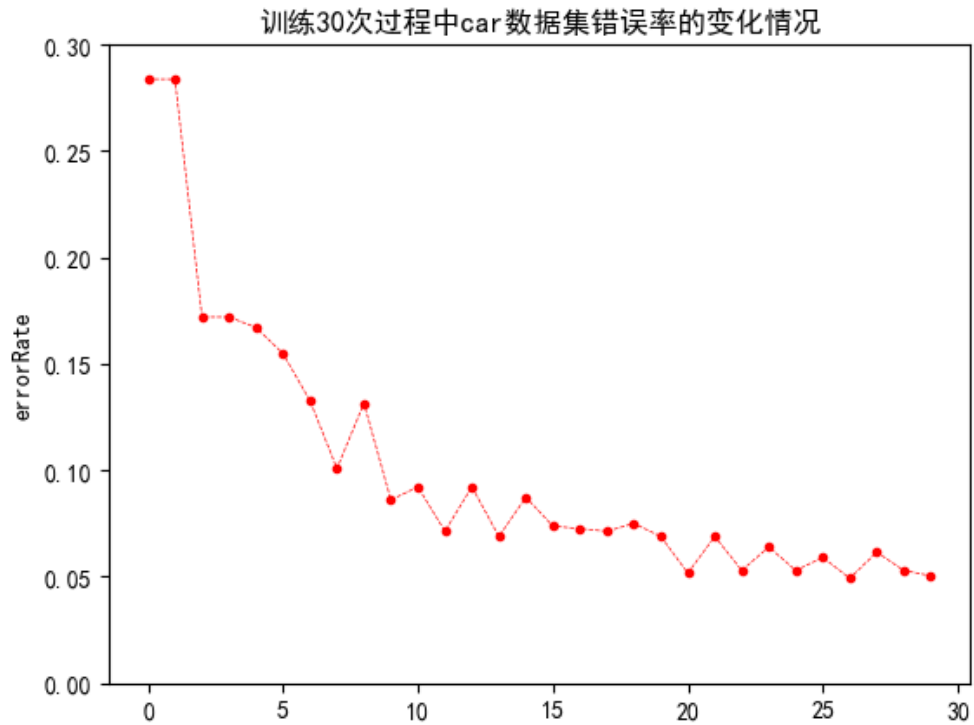


图 5: car30 次训练情况图

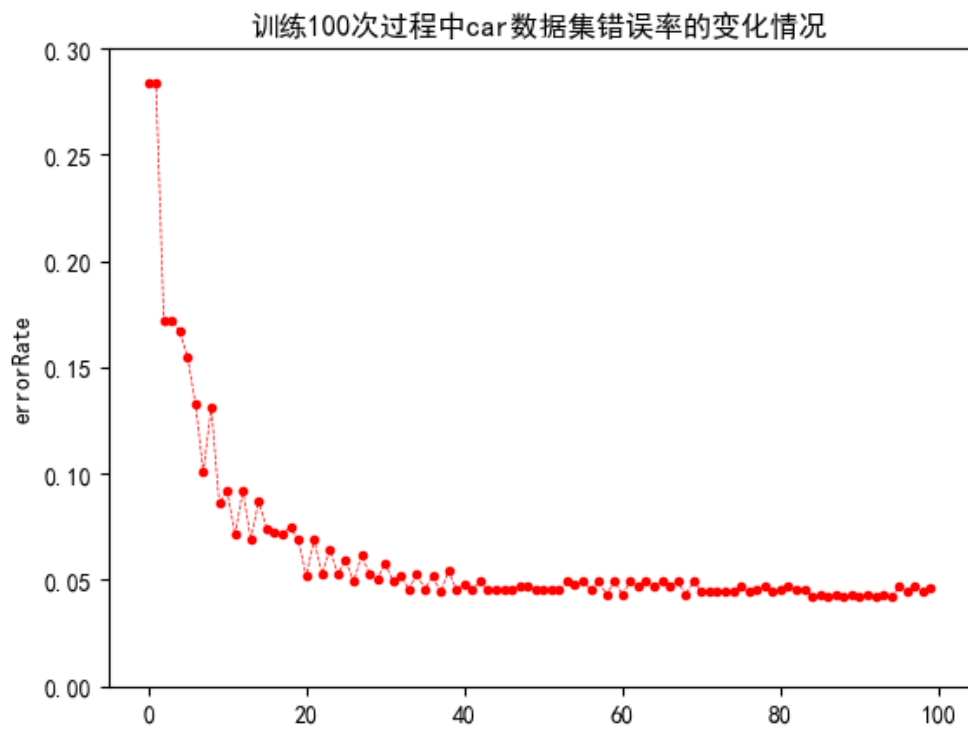


图 6: car100 次训练情况图

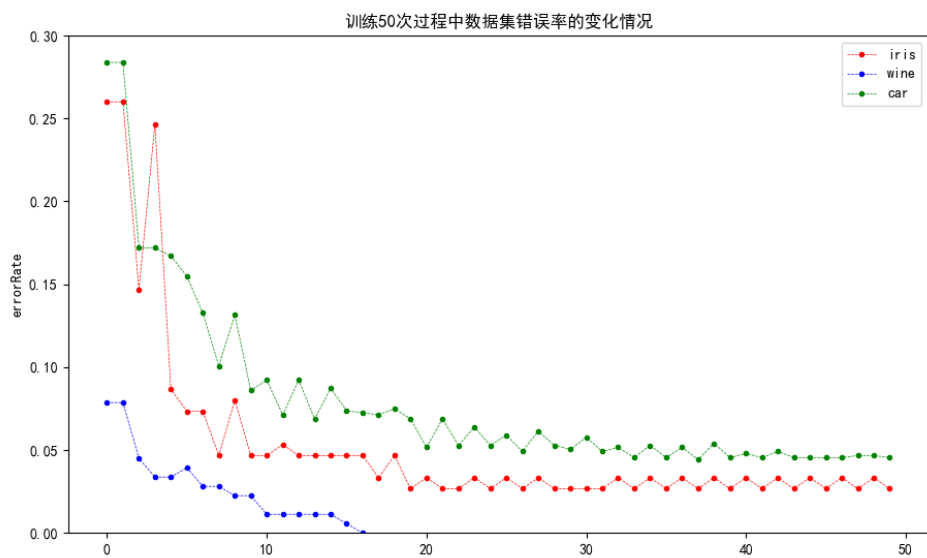


图 7: 三类数据训练情况对比图

第五节 评价方法

十折交叉验证。将数据集分成十份，轮流将其中 9 份作为训练数据，1 份作为测试数据，进行试验。每次试验都会得出相应的正确率（或差错率）。10 次的结果的正确率（或差错率）的平均值作为对算法精度的估计，一般还需要进行多次 10 折交叉验证（例如 10 次 10 折交叉验证），再求其均值，作为对算法准确性的估计。

第六节 实验分析和比较

对于鸢尾花数据进行 10 次十折交叉验证之后的结果为：0.07333；对于红酒数据结果为：0.9；对于汽车数据结果为：0.08588。可以看出此算法相比较而言，可能对于汽车数据更为合适。另外，这次 adaboost 是实现了其二分类的功能，所以在使用的时候需要对数据进行合并，本次采取的合并方式是选择类标号中的一个类作为 1，其他作为-1。进行多次对比实验，寻找训练的 error 结果最佳的。结果图片中对此有所展示。在学习算法的过程中，对过程的理解一开始花了很长的时间，细节上有很多容易搞混，比如其实 adaboost 算法在调权的时候其实调整了两个权重，一个权重是对于样本而言的，一般以权重向量的方式出现，另外一个对于弱分类器而言的。对于分类器的形成方式有很多，此次实现的方式是对于训练集的不同属性的子集训练而成的基分类器。