

摘要

前面的文章已经介绍了五种不同的分类器，它们各有优缺点。我们可以很自然地将不同的分类器组合起来，而这种组合结果则被成为集成方法(ensemble method)或者元算法(meta-algorithm)。使用集成方法时会有多种形式：可以是不同算法的集成，也可以是同一种算法在不同设置下的集成，还可以是数据集不同部分分配给不同分类器之后的集成。



一、前言

前面的文章已经介绍了五种不同的分类器，它们各有优缺点。我们可以很自然地将不同的分类器组合起来，而这种组合结果则被成为**集成方法(ensemble method)**或者**元算法(meta-algorithm)**。使用集成方法时会有多种形式：可以是不同算法的集成，也可以是同一种算法在不同设置下的集成，还可以是数据集不同部分分配给不同分类器之后的集成。

本文出现的所有代码和数据集，均可在我的github上下载，欢迎Follow、Star：[点我查看](#)

二、集成方法

集成方法 (ensemble method) 通过组合多个基分类器 (base classifier) 来完成学习任务，颇有点“三个臭皮匠顶个诸葛亮”的意味。基分类器—弱可学习 (weakly learnable) 分类器，通过集成方法，组合成一个强可学习 (strongly learnable) 分类器。所谓弱可学习，是指学习的正确率仅略优于随机式学习算法；强可学习指正确率较高的多项式学习算法。集成学习的泛化能力一般比单一的基分类器要好，这是因为大部分基分类器都分类错误的概率远离0的。

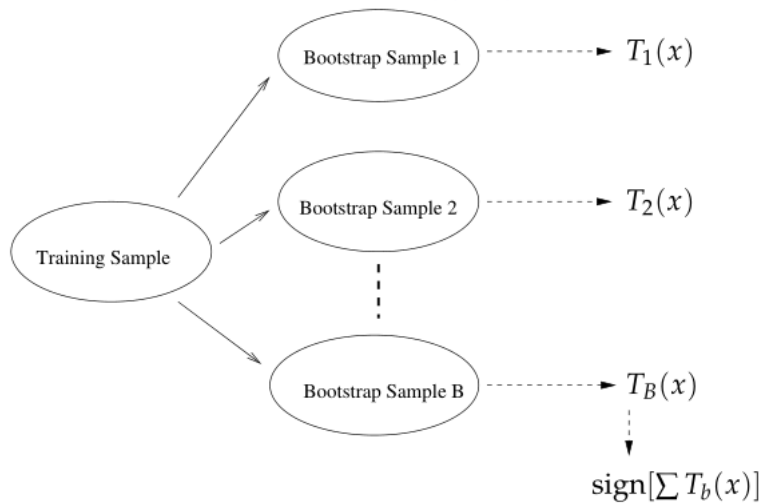
集成方法主要包括Bagging和Boosting两种方法，Bagging和Boosting都是将已有的分类或回归算法通过一定方式组合起来，形成一个性能更加强大准确的。准确的说这是一种分类算法的组装方法，即将弱分类器组装成强分类器的方法。

1、Bagging

自举汇聚法 (bootstrap aggregating)，也称为bagging方法。Bagging对训练数据采用自举采样 (bootstrap sampling)，即有放回地采样数据

- 从原始样本集中抽取训练集。每轮从原始样本集中使用Bootstrapping的方法抽取n个训练样本（在训练集中，有些样本可能被多次抽取到，而有些样本一次都没有被抽中）。共进行k轮抽取，得到k个训练集。（k个训练集之间是相互独立的）
- 每次使用一个训练集得到一个模型，k个训练集共得到k个模型。（注：这里并没有具体的分类算法或回归方法，我们可以根据具体问题采用不同的回归方法，如决策树、感知器等）
- 对分类问题：将上步得到的k个模型采用投票的方式得到分类结果；对回归问题，计算上述模型的均值作为最后的结果。（所有模型的重要性相同）

Schematics of Bagging

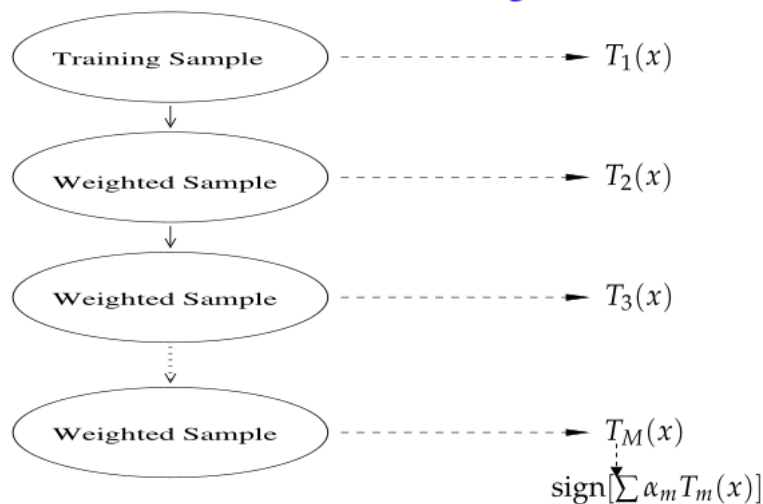


2、Boosting

Boosting是一种与Bagging很类似的技术。Boosting的思路则是采用重赋权（re-weighting）法迭代地训练基分类器，主要思想：

- 每一轮的训练数据样本赋予一个权重，并且每一轮样本的权值分布依赖上一轮的分类结果。
- 基分类器之间采用序列式的线性加权方式进行组合。

Schematics of Boosting



3、Bagging、Boosting二者之间的区别

样本选择上：

- Bagging：训练集是在原始集中有放回选取的，从原始集中选出的各轮训练集之间是独立的。
- Boosting：每一轮的训练集不变，只是训练集中每个样例在分类器中的权重发生变化。而权值是根据上一轮的分类结果进行调整。

样例权重：

- Bagging：使用均匀取样，每个样例的权重相等。
- Boosting：根据错误率不断调整样例的权值，错误率越大则权重越大。

预测函数：

- Bagging：所有预测函数的权重相等。
- Boosting：每个弱分类器都有相应的权重，对于分类误差小的分类器会有更大的权重。

并行计算：

- Bagging：各个预测函数可以并行生成。
- Boosting：各个预测函数只能顺序生成，因为后一个模型参数需要前一轮模型的结果。

4、总结

这两种方法都是把若干个分类器整合为一个分类器的方法，只是整合的方式不一样，最终得到不一样的效果，将不同的分类算法套入到此类算法框架，提高了原单一分类器的分类效果，但是也增大了计算量。

下面是将决策树与这些算法框架进行结合所得到的新的算法：

- Bagging + 决策树 = 随机森林
- AdaBoost + 决策树 = 提升树
- Gradient Boosting + 决策树 = GBDT

集成方法众多，本文主要关注Boosting方法中的一种最流行的版本，即AdaBoost。

三、AdaBoost

AdaBoost算法是基于Boosting思想的机器学习算法，AdaBoost是adaptive boosting（自适应boosting）的缩写，其运行过程如下：

1、计算样本权重

训练数据中的每个样本，赋予其权重，即样本权重，用向量D表示，这些权重都初始化成相等值。假设有n个样本的训练集：

$$\{(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)\}$$

设定每个样本的权重都是相等的，即 $1/n$ 。

2、计算错误率

利用第一个弱学习算法 h_1 对其进行学习，学习完成后进行错误率 ϵ 的统计：

$$\epsilon = \frac{\text{未正确分类的样本数目}}{\text{所有样本数目}}$$

3、计算弱学习算法权重

弱学习算法也有一个权重，用向量 α 表示，利用错误率计算权重 α ：

$$\alpha = \frac{1}{2} \ln\left(\frac{1 - \epsilon}{\epsilon}\right)$$

4、更新样本权重

在第一次学习完成后，需要重新调整样本的权重，以使得在第一分类中被错分的样本的权重，在接下来的学习中可以重点对其进行学习：

$$D_{t+1}(i) = \frac{D_t(i)}{Z_t} \times \begin{cases} e^{-\alpha t} & \text{if } h_t(x_i) = y_i \\ e^{\alpha t} & \text{if } h_t(x_i) \neq y_i \end{cases}$$

其中， $h_t(x_i) = y_i$ 表示对第 i 个样本训练正确，不等于则表示分类错误。 Z_t 是一个归一化因子：

$$Z_t = \text{sum}(D)$$

这个公式我们可以继续化简，将两个公式进行合并，化简如下：

$$D_{t+1}(i) = \frac{D_t(i) \exp(-\alpha_t y_i h_t(x_i))}{\text{sum}(D)}$$

5、AdaBoost算法

重复进行学习，这样经过 t 轮的学习后，就会得到 t 个弱学习算法、权重、弱分类器的输出以及最终的AdaBoost算法的输出，分别如下：

$$\begin{aligned} &\{h_1, h_2, \dots, h_t\} \\ &\{\alpha_1, \alpha_2, \dots, \alpha_t\} \\ &\{h_1(X), h_2(X), \dots, h_t(X)\} \\ &H(X) = \text{sign}\left(\sum_{i=1}^t \alpha_i h_i(X)\right) \end{aligned}$$

其中， $\text{sign}(x)$ 是符号函数。具体过程如下所示：

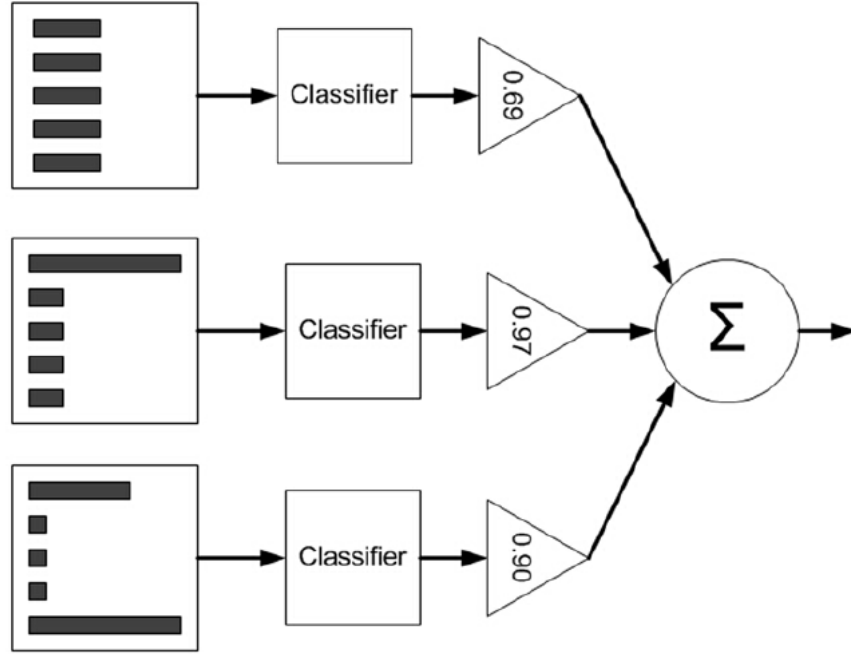


Figure 7.1 Schematic representation of AdaBoost; with the dataset on the left side, the different widths of the bars represent weights applied to each instance. The weighted predictions pass through a classifier, which is then weighted by the triangles (α values). The weighted output of each triangle is summed up in the circle, which produces the final output.

AdaBoost算法总结如下：

Given: $(x_1, y_1), \dots, (x_m, y_m)$ where $x_i \in X, y_i \in Y = \{-1, +1\}$

Initialize $D_1(i) = 1/m$.

For $t = 1, \dots, T$:

- Train weak learner using distribution D_t .
- Get weak hypothesis $h_t : X \rightarrow \{-1, +1\}$ with error

$$\epsilon_t = \Pr_{i \sim D_t} [h_t(x_i) \neq y_i].$$

- Choose $\alpha_t = \frac{1}{2} \ln \left(\frac{1 - \epsilon_t}{\epsilon_t} \right)$.
- Update:

$$\begin{aligned} D_{t+1}(i) &= \frac{D_t(i)}{Z_t} \times \begin{cases} e^{-\alpha_t} & \text{if } h_t(x_i) = y_i \\ e^{\alpha_t} & \text{if } h_t(x_i) \neq y_i \end{cases} \\ &= \frac{D_t(i) \exp(-\alpha_t y_i h_t(x_i))}{Z_t} \end{aligned}$$

where Z_t is a normalization factor (chosen so that D_{t+1} will be a distribution).

Output the final hypothesis:

$$H(x) = \text{sign} \left(\sum_{t=1}^T \alpha_t h_t(x) \right).$$

Figure 1: The boosting algorithm AdaBoost.

四、基于单层决策树构建弱分类器

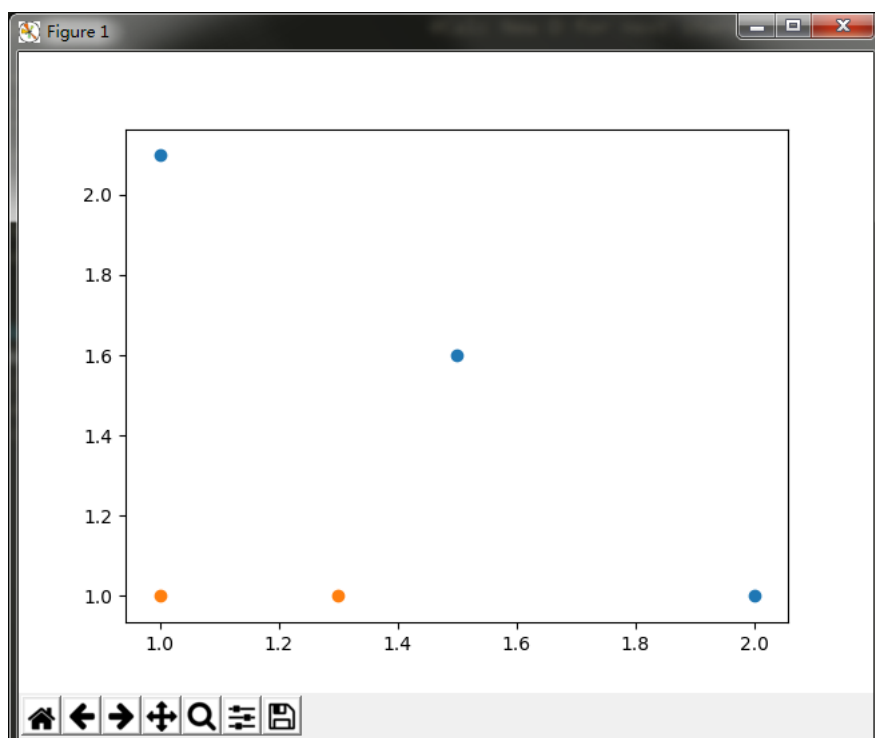
建立AdaBoost算法之前，我们必须先建立弱分类器，并保存样本的权重。弱分类器使用单层决策树（decision stump），也称决策树桩，它是一种通过给定的阈值，进行分类。

1、数据集可视化

为了训练单层决策树，我们需要创建一个训练集，编写代码如下：

```
1 # -*-coding:utf-8 -*-
2 import numpy as np
3 import matplotlib.pyplot as plt
4 """
5 Author:
6     Jack Cui
7 Blog:
8     http://blog.csdn.net/c406495762
9 Zhihu:
10    https://www.zhihu.com/people/Jack--Cui/
11 Modify:
12    2017-10-10
13 """
14 def loadSimpData():
15     """
16     创建单层决策树的数据集
17     Parameters:
18         无
19     Returns:
20         dataMat - 数据矩阵
21         classLabels - 数据标签
22     """
23     dataMat = np.matrix([[ 1. , 2.1],
24                          [ 1.5, 1.6],
25                          [ 1.3, 1. ],
26                          [ 1. , 1. ],
27                          [ 2. , 1. ]])
28     classLabels = [1.0, 1.0, -1.0, -1.0, 1.0]
29     return dataMat, classLabels
30 def showDataSet(dataMat, labelMat):
31     """
32     数据可视化
33     Parameters:
34         dataMat - 数据矩阵
35         labelMat - 数据标签
36     Returns:
37         无
38     """
39     data_plus = []
40     data_minus = []
41     for i in range(len(dataMat)):
42         if labelMat[i] > 0:
43             data_plus.append(dataMat[i])
44         else:
45             data_minus.append(dataMat[i])
46     data_plus_np = np.array(data_plus)
47     data_minus_np = np.array(data_minus)
48     plt.scatter(np.transpose(data_plus_np)[0], np.transpose(data_plus_np)[1])
49     plt.scatter(np.transpose(data_minus_np)[0], np.transpose(data_minus_np)[1])
50     plt.show()
51
52 if __name__ == '__main__':
53     dataArr, classLabels = loadSimpData()
54     showDataSet(dataArr, classLabels)
```

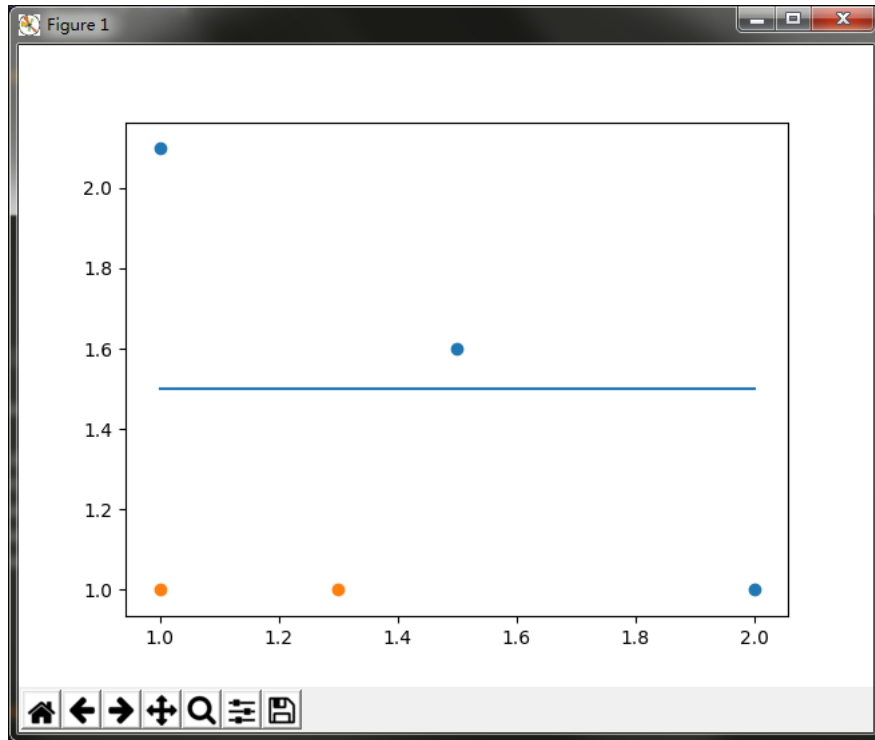
代码运行结果如下：



可以看到，如果想要试着从某个坐标轴上选择一个值（即选择一条与坐标轴平行的直线）来将所有的蓝色圆点和橘色圆点分开，这显然是不可能的。树难以处理的一个著名问题。通过使用多颗单层决策树，我们可以构建出一个能够对该数据集完全正确分类的分类器。

2、构建单层决策树

我们设置一个分类阈值，比如我横向切分，如下图所示：



蓝横线上边的是一个类别，蓝横线下边是一个类别。显然，此时有一个蓝点分类错误，计算此时的分类误差，误差为 $1/5 = 0.2$ 。这个横线与坐标轴的是我们设置的阈值，通过不断改变阈值的大小，找到使单层决策树的分类误差最小的阈值。同理，竖线也是如此，找到最佳分类的阈值，就找到了最佳单代码如下：

```

1  #-*-coding:utf-8 -*-
2  import numpy as np
3  import matplotlib.pyplot as plt
4  """
5  Author:
6      Jack Cui
7  Blog:
8      http://blog.csdn.net/c406495762
9  Zhihu:
10     https://www.zhihu.com/people/Jack--Cui/
11  Modify:
12     2017-10-10
13  """
14  def loadSimpData():
15      """
16      创建单层决策树的数据集
17      Parameters:
18          无
19      Returns:
20          dataMat - 数据矩阵
21          classLabels - 数据标签
22      """
23      datMat = np.matrix([[ 1. ,  2.1],
24                          [ 1.5,  1.6],
25                          [ 1.3,  1. ],
26                          [ 1. ,  1. ],
27                          [ 2. ,  1. ]])
28      classLabels = [1.0, 1.0, -1.0, -1.0, 1.0]
29      return datMat, classLabels
30
31  def stumpClassify(dataMatrix, dimen, threshVal, threshIneq):
32      """
33      单层决策树分类函数
34      Parameters:
35          dataMatrix - 数据矩阵
36          dimen - 第dimen列，也就是第几个特征
37          threshVal - 阈值
38          threshIneq - 标志
39      Returns:
40          retArray - 分类结果
41      """
42      retArray = np.ones((np.shape(dataMatrix)[0],1)) #初始化retArray为1
43      if threshIneq == 'lt':
44          retArray[dataMatrix[:,dimen] <= threshVal] = -1.0 #如果小于阈值，则赋值为-1
45      else:
46          retArray[dataMatrix[:,dimen] > threshVal] = -1.0 #如果大于阈值，则赋值为-1
47      return retArray
48
49  def buildStump(dataArr, classLabels, D):

```

```

50     """
51     找到数据集上最佳的单层决策树
52     Parameters:
53         dataArr - 数据矩阵
54         classLabels - 数据标签
55         D - 样本权重
56     Returns:
57         bestStump - 最佳单层决策树信息
58         minError - 最小误差
59         bestClasEst - 最佳的分类结果
60     """
61     dataMatrix = np.mat(dataArr); labelMat = np.mat(classLabels).T
62     m,n = np.shape(dataMatrix)
63     numSteps = 10.0; bestStump = {}; bestClasEst = np.mat(np.zeros((m,1)))
64     minError = float('inf')
65     for i in range(n):
66         rangeMin = dataMatrix[:,i].min(); rangeMax = dataMatrix[:,i].max()
67         stepSize = (rangeMax - rangeMin) / numSteps
68         for j in range(-1, int(numSteps) + 1):
69             for inequal in ['lt', 'gt']:
70                 threshVal = (rangeMin + float(j) * stepSize)
71                 predictedVals = stumpClassify(dataMatrix, i, threshVal, inequal)
72                 errArr = np.mat(np.ones((m,1)))
73                 errArr[predictedVals == labelMat] = 0
74                 weightedError = D.T * errArr
75                 print("split: dim %d, thresh %.2f, thresh inequal: %s, the weighted error is %.3f" % (i, threshVal, inequal, weightedError))
76                 if weightedError < minError:
77                     minError = weightedError
78                     bestClasEst = predictedVals.copy()
79                     bestStump['dim'] = i
80                     bestStump['thresh'] = threshVal
81                     bestStump['ineq'] = inequal
82     return bestStump,minError,bestClasEst
83
84 if __name__ == '__main__':
85     dataArr,classLabels = loadSimpData()
86     D = np.mat(np.ones((5, 1)) / 5)
87     bestStump,minError,bestClasEst = buildStump(dataArr,classLabels,D)
88     print('bestStump:\n', bestStump)
89     print('minError:\n', minError)
90     print('bestClasEst:\n', bestClasEst)

```

代码运行结果如下：

```

133 if __name__ == '__main__':
134     dataArr,classLabels = loadSimpData()
135     D = np.mat(np.ones((5, 1)) / 5)
136     bestStump,minError,bestClasEst = buildStump(dataArr,classLabels,D)
137     print('bestStump:\n', bestStump)
138     print('minError:\n', minError)
139     print('bestClasEst:\n', bestClasEst)

split: dim 1, thresh 1.44, thresh inequal: gt, the weighted error is 0.800
split: dim 1, thresh 1.55, thresh inequal: lt, the weighted error is 0.200
split: dim 1, thresh 1.55, thresh inequal: gt, the weighted error is 0.800
split: dim 1, thresh 1.66, thresh inequal: lt, the weighted error is 0.400
split: dim 1, thresh 1.66, thresh inequal: gt, the weighted error is 0.600
split: dim 1, thresh 1.77, thresh inequal: lt, the weighted error is 0.400
split: dim 1, thresh 1.77, thresh inequal: gt, the weighted error is 0.600
split: dim 1, thresh 1.88, thresh inequal: lt, the weighted error is 0.400
split: dim 1, thresh 1.88, thresh inequal: gt, the weighted error is 0.600
split: dim 1, thresh 1.99, thresh inequal: lt, the weighted error is 0.400
split: dim 1, thresh 1.99, thresh inequal: gt, the weighted error is 0.600
split: dim 1, thresh 2.10, thresh inequal: lt, the weighted error is 0.600
split: dim 1, thresh 2.10, thresh inequal: gt, the weighted error is 0.400
bestStump:
{'thresh': 1.3, 'ineq': 'lt', 'dim': 0}
minError:
[[ 0.2]]
bestClasEst:
[[-1.]
 [ 1.]
 [-1.]
 [-1.]
 [ 1.]]
[Finished in 0.6s]

```

代码不难理解，就是通过遍历，改变不同的阈值，计算最终的分类误差，找到分类误差最小的分类方式，即为我们要找的最佳单层决策树。这里lt表示分类方式，对于小于阈值的样本点赋值为-1，gt表示greater than，也是表示分类方式，对于大于阈值的样本点赋值为-1。经过遍历，我们找到，训练策树的最小分类误差为0.2，就是对于该数据集，无论用什么样的单层决策树，分类误差最小就是0.2。这就是我们训练好的弱分类器。接下来，使用AdaB分类器性能，将分类误差缩短到0，看下AdaBoost算法是如何实现的。

五、使用AdaBoost提升分类器性能

根据之前介绍的AdaBoost算法实现过程，使用AdaBoost算法提升分类器性能，编写代码如下：

```

1 # -*-coding:utf-8 -*-
2 import numpy as np
3 import matplotlib.pyplot as plt
4 """

```



```

5 Author:
6     Jack Cui
7 Blog:
8     http://blog.csdn.net/c406495762
9 Zhihu:
10    https://www.zhihu.com/people/Jack--Cui/
11 Modify:
12    2017-10-10
13 """
14 def loadSimpData():
15     """
16     创建单层决策树的数据集
17     Parameters:
18         无
19     Returns:
20         dataMat - 数据矩阵
21         classLabels - 数据标签
22     """
23     datMat = np.matrix([[ 1. , 2.1],
24                          [ 1.5, 1.6],
25                          [ 1.3, 1. ],
26                          [ 1. , 1. ],
27                          [ 2. , 1. ]])
28     classLabels = [1.0, 1.0, -1.0, -1.0, 1.0]
29     return datMat, classLabels
30
31 def stumpClassify(dataMatrix, dimen, threshVal, threshIneq):
32     """
33     单层决策树分类函数
34     Parameters:
35         dataMatrix - 数据矩阵
36         dimen - 第dimen列，也就是第几个特征
37         threshVal - 阈值
38         threshIneq - 标志
39     Returns:
40         retArray - 分类结果
41     """
42     retArray = np.ones((np.shape(dataMatrix)[0], 1)) #初始化retArray为1
43     if threshIneq == 'lt':
44         retArray[dataMatrix[:, dimen] <= threshVal] = -1.0 #如果小于阈值，则赋值为-1
45     else:
46         retArray[dataMatrix[:, dimen] > threshVal] = -1.0 #如果大于阈值，则赋值为-1
47     return retArray
48
49 def buildStump(dataArr, classLabels, D):
50     """
51     找到数据集上最佳的单层决策树
52     Parameters:
53         dataArr - 数据矩阵
54         classLabels - 数据标签
55         D - 样本权重
56     Returns:
57         bestStump - 最佳单层决策树信息
58         minError - 最小误差
59         bestClasEst - 最佳的分类结果
60     """
61     dataMatrix = np.mat(dataArr); labelMat = np.mat(classLabels).T
62     m, n = np.shape(dataMatrix)
63     numSteps = 10.0; bestStump = {}; bestClasEst = np.mat(np.zeros((m, 1)))
64     minError = float('inf') #最小误差初始化为正无穷大
65     for i in range(n): #遍历所有特征
66         rangeMin = dataMatrix[:, i].min(); rangeMax = dataMatrix[:, i].max() #找到特征中最小的值和最大值
67         stepSize = (rangeMax - rangeMin) / numSteps #计算步长
68         for j in range(-1, int(numSteps) + 1):
69             for inequal in ['lt', 'gt']: #大于和小于的情况，均遍历。lt:less than, gt:greater than
70                 threshVal = (rangeMin + float(j) * stepSize) #计算阈值
71                 predictedVals = stumpClassify(dataMatrix, i, threshVal, inequal) #计算分类结果
72                 errArr = np.mat(np.ones((m, 1))) #初始化误差矩阵
73                 errArr[predictedVals == labelMat] = 0 #分类正确的，赋值为0
74                 weightedError = D.T * errArr #计算误差
75                 print("split: dim %d, thresh %.2f, thresh inequal: %s, the weighted error is %.3f" % (i, threshVal, inequal, weightedError)) #找到误差最小的分类方式
76                 if weightedError < minError:
77                     minError = weightedError
78                     bestClasEst = predictedVals.copy()
79                     bestStump['dim'] = i
80                     bestStump['thresh'] = threshVal
81                     bestStump['ineq'] = inequal
82     return bestStump, minError, bestClasEst
83
84 def adaBoostTrainDS(dataArr, classLabels, numIt = 40):
85     weakClassArr = []
86     m = np.shape(dataArr)[0]
87     D = np.mat(np.ones((m, 1)) / m) #初始化权重
88     aggClassEst = np.mat(np.zeros((m, 1)))
89     for i in range(numIt):
90         bestStump, error, classEst = buildStump(dataArr, classLabels, D) #构建单层决策树
91         print("D:", D.T)
92         alpha = float(0.5 * np.log((1.0 - error) / max(error, 1e-16))) #计算弱学习算法权重alpha,使error不等于0,因为分母不能为0
93         bestStump['alpha'] = alpha #存储弱学习算法权重
94         weakClassArr.append(bestStump) #存储单层决策树
95         print("classEst: ", classEst.T)
96         expon = np.multiply(-1 * alpha * np.mat(classLabels).T, classEst) #计算e的指数项
97         D = np.multiply(D, np.exp(expon))
98         D = D / D.sum() #根据样本权重公式，更新样本权重
99         #计算AdaBoost误差，当误差为0的时候，退出循环
100        aggClassEst += alpha * classEst
101        print("aggClassEst: ", aggClassEst.T)
102        aggErrors = np.multiply(np.sign(aggClassEst) != np.mat(classLabels).T, np.ones((m, 1))) #计算误差
103        errorRate = aggErrors.sum() / m
104        print("total error: ", errorRate)
105        if errorRate == 0.0: break #误差为0，退出循环
106    return weakClassArr, aggClassEst
107
108 if __name__ == '__main__':

```



```

109 dataArr,classLabels = loadSimpData()
110 weakClassArr, aggClassEst = adaBoostTrainDS(dataArr, classLabels)
111 print(weakClassArr)
112 print(aggClassEst)

```

运行结果如下：

```

split: dim 1, thresh 1.11, thresh inequal: gt, the weighted error is 0.500
split: dim 1, thresh 1.22, thresh inequal: lt, the weighted error is 0.500
split: dim 1, thresh 1.22, thresh inequal: gt, the weighted error is 0.500
split: dim 1, thresh 1.33, thresh inequal: lt, the weighted error is 0.500
split: dim 1, thresh 1.33, thresh inequal: gt, the weighted error is 0.500
split: dim 1, thresh 1.44, thresh inequal: lt, the weighted error is 0.500
split: dim 1, thresh 1.44, thresh inequal: gt, the weighted error is 0.500
split: dim 1, thresh 1.55, thresh inequal: lt, the weighted error is 0.500
split: dim 1, thresh 1.55, thresh inequal: gt, the weighted error is 0.500
split: dim 1, thresh 1.66, thresh inequal: lt, the weighted error is 0.571
split: dim 1, thresh 1.66, thresh inequal: gt, the weighted error is 0.429
split: dim 1, thresh 1.77, thresh inequal: lt, the weighted error is 0.571
split: dim 1, thresh 1.77, thresh inequal: gt, the weighted error is 0.429
split: dim 1, thresh 1.88, thresh inequal: lt, the weighted error is 0.571
split: dim 1, thresh 1.88, thresh inequal: gt, the weighted error is 0.429
split: dim 1, thresh 1.99, thresh inequal: lt, the weighted error is 0.571
split: dim 1, thresh 1.99, thresh inequal: gt, the weighted error is 0.429
split: dim 1, thresh 2.10, thresh inequal: lt, the weighted error is 0.857
split: dim 1, thresh 2.10, thresh inequal: gt, the weighted error is 0.143
D: [[ 0.28571429  0.07142857  0.07142857  0.07142857  0.5        ]]
classEst: [[ 1.  1.  1.  1.  1.]]
aggClassEst: [[ 1.17568763  2.56198199 -0.77022252 -0.77022252  0.61607184]]
total error: 0.0
[{'ineq': 'lt', 'dim': 0, 'alpha': 0.6931471805599453, 'thresh': 1.3}, {'ineq': 'lt', 'dim':
[ 1.17568763]
[ 2.56198199]
[-0.77022252]
[-0.77022252]
[ 0.61607184]]

```

在第一轮迭代中，D中的所有值都相等。于是，只有第一个数据点被错分了。因此在第二轮迭代中，D向量给第一个数据点0.5的权重。这就可以通过aggClassEst的符号来了解总的类别。第二次迭代之后，我们会发现第一个数据点已经正确分类了，但此时最后一个数据点却是错分了。D向量中的最后0.5，而D向量中的其他值都变得非常小。最后，第三次迭代之后aggClassEst所有值的符号和真是类别标签都完全吻合，那么训练错误率为0，程序终止运行。

最后训练结果包含了三个弱分类器，其中包含了分类所需要的所有信息。一共迭代了3次，所以训练了3个弱分类器构成一个使用AdaBoost算法优化分类器的错误率为0。

一旦拥有了多个弱分类器以及其对应的alpha值，进行测试就变得相当容易了。编写代码如下：

```

1 # -*-coding:utf-8 -*-
2 import numpy as np
3 import matplotlib.pyplot as plt
4 """
5 Author:
6 Jack Cui
7 Blog:
8 http://blog.csdn.net/c406495762
9 Zhihu:
10 https://www.zhihu.com/people/Jack--Cui/
11 Modify:
12 2017-10-10
13 """
14 def loadSimpData():
15     """
16     创建单层决策树的数据集
17     Parameters:
18     无
19     Returns:
20     dataMat - 数据矩阵
21     classLabels - 数据标签
22     """
23     dataMat = np.matrix([[ 1. ,  2.1],
24                          [ 1.5,  1.6],
25                          [ 1.3,  1. ],
26                          [ 1. ,  1. ],
27                          [ 2. ,  1. ]])
28     classLabels = [1.0, 1.0, -1.0, -1.0, 1.0]
29     return dataMat, classLabels
30 def showDataSet(dataMat, labelMat):
31     """
32     数据可视化
33     Parameters:
34     dataMat - 数据矩阵
35     labelMat - 数据标签
36     Returns:
37     无
38     """
39     data_plus = [] #正样本
40     data_minus = [] #负样本
41     for i in range(len(dataMat)):
42         if labelMat[i] > 0:
43             data_plus.append(dataMat[i])
44         else:
45             data_minus.append(dataMat[i])

```

```

46 data_plus_np = np.array(data_plus) #转换为numpy矩阵
47 data_minus_np = np.array(data_minus) #转换为numpy矩阵
48 plt.scatter(np.transpose(data_plus_np)[0], np.transpose(data_plus_np)[1]) #正样本散点图
49 plt.scatter(np.transpose(data_minus_np)[0], np.transpose(data_minus_np)[1]) #负样本散点图
50 plt.show()
51 def stumpClassify(dataMatrix, dimen, threshVal, threshIneq):
52     """
53     单层决策树分类函数
54     Parameters:
55         dataMatrix - 数据矩阵
56         dimen - 第dimen列, 也就是第几个特征
57         threshVal - 阈值
58         threshIneq - 标志
59     Returns:
60         retArray - 分类结果
61     """
62     retArray = np.ones((np.shape(dataMatrix)[0],1)) #初始化retArray为1
63     if threshIneq == 'lt':
64         retArray[dataMatrix[:,dimen] <= threshVal] = -1.0 #如果小于阈值,则赋值为-1
65     else:
66         retArray[dataMatrix[:,dimen] > threshVal] = -1.0 #如果大于阈值,则赋值为-1
67     return retArray
68
69 def buildStump(dataArr, classLabels, D):
70     """
71     找到数据集上最佳的单层决策树
72     Parameters:
73         dataArr - 数据矩阵
74         classLabels - 数据标签
75         D - 样本权重
76     Returns:
77         bestStump - 最佳单层决策树信息
78         minError - 最小误差
79         bestClasEst - 最佳的分类结果
80     """
81     dataMatrix = np.mat(dataArr); labelMat = np.mat(classLabels).T
82     m, n = np.shape(dataMatrix)
83     numSteps = 10.0; bestStump = {}; bestClasEst = np.mat(np.zeros((m,1)))
84     minError = float('inf') #最小误差初始化为正无穷大
85     for i in range(n): #遍历所有特征
86         rangeMin = dataMatrix[:,i].min(); rangeMax = dataMatrix[:,i].max() #找到特征中最小的值和最大值
87         stepSize = (rangeMax - rangeMin) / numSteps #计算步长
88         for j in range(-1, int(numSteps) + 1):
89             for inequal in ['lt', 'gt']: #大于和小于的情况, 均遍历。lt:less than, gt:greater than
90                 threshVal = (rangeMin + float(j) * stepSize) #计算阈值
91                 predictedVals = stumpClassify(dataMatrix, i, threshVal, inequal) #计算分类结果
92                 errArr = np.mat(np.ones((m,1))) #初始化误差矩阵
93                 errArr[predictedVals == labelMat] = 0 #分类正确的,赋值为0
94                 weightedError = D.T * errArr #计算误差
95                 # print("split: dim %d, thresh %.2f, thresh inequal: %s, the weighted error is %.3f" % (i, threshVal, inequal, weightedError))
96                 if weightedError < minError: #找到误差最小的分类方式
97                     minError = weightedError
98                     bestClasEst = predictedVals.copy()
99                     bestStump['dim'] = i
100                     bestStump['thresh'] = threshVal
101                     bestStump['ineq'] = inequal
102     return bestStump, minError, bestClasEst
103 def adaBoostTrainDS(dataArr, classLabels, numIt = 40):
104     """
105     使用AdaBoost算法提升弱分类器性能
106     Parameters:
107         dataArr - 数据矩阵
108         classLabels - 数据标签
109         numIt - 最大迭代次数
110     Returns:
111         weakClassArr - 训练好的分类器
112         aggClassEst - 类别估计累计值
113     """
114     weakClassArr = []
115     m = np.shape(dataArr)[0]
116     D = np.mat(np.ones((m, 1)) / m) #初始化权重
117     aggClassEst = np.mat(np.zeros((m,1)))
118     for i in range(numIt):
119         bestStump, error, classEst = buildStump(dataArr, classLabels, D) #构建单层决策树
120         # print("D:",D.T)
121         alpha = float(0.5 * np.log((1.0 - error) / max(error, 1e-16))) #计算弱学习算法权重alpha,使error不等于0,因为分母不能为0
122         bestStump['alpha'] = alpha #存储弱学习算法权重
123         weakClassArr.append(bestStump) #存储单层决策树
124         # print("classEst: ", classEst.T)
125         expon = np.multiply(-1 * alpha * np.mat(classLabels).T, classEst) #计算e的指数项
126         D = np.multiply(D, np.exp(expon))
127         D = D / D.sum() #根据样本权重公式,更新样本权重
128         #计算AdaBoost误差. 当误差为0的时候,退出循环
129         aggClassEst += alpha * classEst #计算类别估计累计值
130         # print("aggClassEst: ", aggClassEst.T)
131         aggErrors = np.multiply(np.sign(aggClassEst) != np.mat(classLabels).T, np.ones((m,1))) #计算误差
132         errorRate = aggErrors.sum() / m
133         # print("total error: ", errorRate)
134         if errorRate == 0.0: break #误差为0,退出循环
135     return weakClassArr, aggClassEst
136 def adaClassify(datToClass, classifierArr):
137     """
138     AdaBoost分类函数
139     Parameters:
140         datToClass - 待分类样例
141         classifierArr - 训练好的分类器
142     Returns:
143         分类结果
144     """
145     dataMatrix = np.mat(datToClass)
146     m = np.shape(dataMatrix)[0]
147     aggClassEst = np.mat(np.zeros((m,1)))
148     for i in range(len(classifierArr)): #遍历所有分类器,进行分类
149         classEst = stumpClassify(dataMatrix, classifierArr[i]['dim'], classifierArr[i]['thresh'], classifierArr[i]['ineq'])

```

```

150     aggClassEst += classifierArr[i]['alpha'] * classEst
151     print(aggClassEst)
152     return np.sign(aggClassEst)
153 if __name__ == '__main__':
154     dataArr, classLabels = loadSimpData()
155     weakClassArr, aggClassEst = adaBoostTrainDS(dataArr, classLabels)
156     print(adaClassify([[0,0],[5,5]], weakClassArr))

```

运行结果如下图所示：

```

162 if __name__ == '__main__':
163     dataArr, classLabels = loadSimpData()
164     weakClassArr, aggClassEst = adaBoostTrainDS(dataArr, classLabels)
165     print(adaClassify([[0,0],[5,5]], weakClassArr))

```

```

[[-0.69314718]
 [ 0.69314718]]
[[-1.66610226]
 [ 1.66610226]]
[[-2.56198199]
 [ 2.56198199]]
[[-1.]
 [ 1.]]
[Finished in 0.8s]

```

代码很简单，在之前代码的基础上，添加adaClassify()函数，该函数遍历所有训练得到的弱分类器，利用单层决策树，输出的类别估计值乘以该单层权重alpha，然后累加到aggClassEst上，最后通过sign函数最终的结果。可以看到，分类没有问题，(5,5)属于正类，(0,0)属于负类。

六、在一个难数据集上应用AdaBoost

在《Python3《机器学习实战》学习笔记（七）：Logistic回归实战篇之预测病马死亡率》文章中，我们使用Logistic回归方法训练马疝病数据集，预率。当时的训练结果如下图所示：

```

185     testSet.append(lineArr)
186     testLabels.append(float(currLine[-1]))
187     classifier = LogisticRegression(solver='sag', max_iter=5000).fit(trainingSet, testLabels)
188     test_accuracy = classifier.score(testSet, testLabels) * 100
189     print('正确率:%f%%' % test_accuracy)
190
191 if __name__ == '__main__':
192     colicSklearn()

```

```

正确率:73.134328%
[Finished in 1.4s]

```

这个使用Sklearn的LogisticRegression()训练的分类器，可以看到，正确率约为73.134%，也就是说错误率约为26.866%。可以看到错误率还是蛮高的，使用AdaBoost算法，训练出一个更强的分类器，这里的数据集有所变化，之前的标签是0和1，现在将标签改为+1和-1，其他数据不变。

更改好的数据集下载地址：[数据集下载](#)

1、自己动手丰衣足食

使用自己的用Python写的AdaBoost算法进行训练，添加loadDataSet函数用于加载数据集。编写代码如下：

```

1  # -*-coding:utf-8 -*-
2  import numpy as np
3  import matplotlib.pyplot as plt
4  """
5  Author:
6  Jack Cui
7  Blog:
8  http://blog.csdn.net/c406495762
9  Zhihu:
10 https://www.zhihu.com/people/Jack--Cui/
11 Modify:
12 2017-10-10
13 """
14 def loadDataSet(fileName):
15     numFeat = len((open(fileName).readline().split('\t')))
16     dataMat = []; labelMat = []
17     fr = open(fileName)
18     for line in fr.readlines():
19         lineArr = []
20         curLine = line.strip().split('\t')
21         for i in range(numFeat - 1):
22             lineArr.append(float(curLine[i]))
23         dataMat.append(lineArr)
24         labelMat.append(float(curLine[-1]))
25     return dataMat, labelMat
26 def stumpClassify(dataMatrix, dimen, threshVal, threshIneq):
27     """
28     单层决策树分类函数

```

```

29 Parameters:
30     dataMatrix - 数据矩阵
31     dimen - 第dimen列，也就是第几个特征
32     threshVal - 阈值
33     threshIneq - 标志
34 Returns:
35     retArray - 分类结果
36     """
37     retArray = np.ones((np.shape(dataMatrix)[0],1)) #初始化retArray为1
38     if threshIneq == 'lt':
39         retArray[dataMatrix[:,dimen] <= threshVal] = -1.0 #如果小于阈值,则赋值为-1
40     else:
41         retArray[dataMatrix[:,dimen] > threshVal] = -1.0 #如果大于阈值,则赋值为-1
42     return retArray
43
44 def buildStump(dataArr,classLabels,D):
45     """
46     找到数据集上最佳的单层决策树
47     Parameters:
48         dataArr - 数据矩阵
49         classLabels - 数据标签
50         D - 样本权重
51     Returns:
52         bestStump - 最佳单层决策树信息
53         minError - 最小误差
54         bestClasEst - 最佳的分类结果
55     """
56     dataMatrix = np.mat(dataArr); labelMat = np.mat(classLabels).T
57     m,n = np.shape(dataMatrix)
58     numSteps = 10.0; bestStump = {}; bestClasEst = np.mat(np.zeros((m,1)))
59     minError = float('inf') #最小误差初始化为正无穷大
60     for i in range(n): #遍历所有特征
61         rangeMin = dataMatrix[:,i].min(); rangeMax = dataMatrix[:,i].max() #找到特征中最小的值和最大值
62         stepSize = (rangeMax - rangeMin) / numSteps #计算步长
63         for j in range(-1, int(numSteps) + 1):
64             for inequal in ['lt', 'gt']: #大于和小于的情况，均遍历。lt:less than, gt:greater than
65                 threshVal = (rangeMin + float(j) * stepSize) #计算阈值
66                 predictedVals = stumpClassify(dataMatrix, i, threshVal, inequal) #计算分类结果
67                 errArr = np.mat(np.ones((m,1))) #初始化误差矩阵
68                 errArr[predictedVals == labelMat] = 0 #分类正确的,赋值为0
69                 weightedError = D.T * errArr #计算误差
70                 # print("split: dim %d, thresh %.2f, thresh inequal: %s, the weighted error is %.3f" % (i, threshVal, inequal, weightedError))
71                 if weightedError < minError: #找到误差最小的分类方式
72                     minError = weightedError
73                     bestClasEst = predictedVals.copy()
74                     bestStump['dim'] = i
75                     bestStump['thresh'] = threshVal
76                     bestStump['ineq'] = inequal
77     return bestStump, minError, bestClasEst
78
79 def adaBoostTrainDS(dataArr, classLabels, numIt = 40):
80     """
81     使用AdaBoost算法提升弱分类器性能
82     Parameters:
83         dataArr - 数据矩阵
84         classLabels - 数据标签
85         numIt - 最大迭代次数
86     Returns:
87         weakClassArr - 训练好的分类器
88         aggClassEst - 类别估计累计值
89     """
90     weakClassArr = []
91     m = np.shape(dataArr)[0]
92     D = np.mat(np.ones((m, 1)) / m) #初始化权重
93     aggClassEst = np.mat(np.zeros((m,1)))
94     for i in range(numIt):
95         bestStump, error, classEst = buildStump(dataArr, classLabels, D) #构建单层决策树
96         # print("D:",D.T)
97         alpha = float(0.5 * np.log((1.0 - error) / max(error, 1e-16))) #计算弱学习算法权重alpha,使error不等于0,因为分母不能为0
98         bestStump['alpha'] = alpha #存储弱学习算法权重
99         weakClassArr.append(bestStump) #存储单层决策树
100         # print("classEst: ", classEst.T)
101         expon = np.multiply(-1 * alpha * np.mat(classLabels).T, classEst) #计算e的指数项
102         D = np.multiply(D, np.exp(expon))
103         D = D / D.sum() #根据样本权重公式，更新样本权重
104         #计算AdaBoost误差，当误差为0的时候，退出循环
105         aggClassEst += alpha * classEst #计算类别估计累计值
106         # print("aggClassEst: ", aggClassEst.T)
107         aggErrors = np.multiply(np.sign(aggClassEst) != np.mat(classLabels).T, np.ones((m,1))) #计算误差
108         errorRate = aggErrors.sum() / m
109         # print("total error: ", errorRate)
110         if errorRate == 0.0: break #误差为0，退出循环
111     return weakClassArr, aggClassEst
112
113 def adaClassify(datToClass,classifierArr):
114     """
115     AdaBoost分类函数
116     Parameters:
117         datToClass - 待分类样例
118         classifierArr - 训练好的分类器
119     Returns:
120         分类结果
121     """
122     dataMatrix = np.mat(datToClass)
123     m = np.shape(dataMatrix)[0]
124     aggClassEst = np.mat(np.zeros((m,1)))
125     for i in range(len(classifierArr)): #遍历所有分类器，进行分类
126         classEst = stumpClassify(dataMatrix, classifierArr[i]['dim'], classifierArr[i]['thresh'], classifierArr[i]['ineq'])
127         aggClassEst += classifierArr[i]['alpha'] * classEst
128         # print(aggClassEst)
129     return np.sign(aggClassEst)
130
131 if __name__ == '__main__':
132     dataArr, labelArr = loadDataSet('horseColicTraining2.txt')

```



```

133 weakClassArr, aggClassEst = adaBoostTrainDS(dataArr, LabelArr)
134 testArr, testLabelArr = loadDataSet('horseColicTest2.txt')
135 print(weakClassArr)
136 predictions = adaClassify(dataArr, weakClassArr)
137 errArr = np.mat(np.ones((len(dataArr), 1)))
138 print('训练集的错误率:%.3f%%' % float(errArr[predictions != np.mat(LabelArr).T].sum() / len(dataArr) * 100))
139 predictions = adaClassify(testArr, weakClassArr)
140 errArr = np.mat(np.ones((len(testArr), 1)))
141 print('测试集的错误率:%.3f%%' % float(errArr[predictions != np.mat(testLabelArr).T].sum() / len(testArr) * 100))

```

代码运行结果如下：

```

135
136 if __name__ == '__main__':
137     dataArr, LabelArr = loadDataSet('horseColicTraining2.txt')
138     weakClassArr, aggClassEst = adaBoostTrainDS(dataArr, LabelArr)
139     testArr, testLabelArr = loadDataSet('horseColicTest2.txt')
140     print(weakClassArr)
141     predictions = adaClassify(dataArr, weakClassArr)
142     errArr = np.mat(np.ones((len(dataArr), 1)))
143     print('训练集的错误率:%.3f%%' % float(errArr[predictions != np.mat(LabelArr).T].sum() / len(dataArr) *
144     predictions = adaClassify(testArr, weakClassArr)
145     errArr = np.mat(np.ones((len(testArr), 1)))
146     print('测试集的错误率:%.3f%%' % float(errArr[predictions != np.mat(testLabelArr).T].sum() / len(testAr
147

```

```

[{'ineq': 'gt', 'dim': 9, 'alpha': 0.4616623792657674, 'thresh': 3.0}, {'ineq': 'gt', 'dim': 17, 'alpha': 0.31248245042467115, 'thresh': 52.5}, {'ineq': 'gt', 'dim': 3, 'alpha': 0.28680973201695786, 'thresh': 55.199999999999996}, {'ineq': 'lt', 'dim': 18, 'alpha': 0.23297004638939506, 'thresh': 62.300000000000004}, {'ineq': 'lt', 'dim': 10, 'alpha': 0.19803846151213736, 'thresh': 0.0}, {'ineq': 'lt', 'dim': 5, 'alpha': 0.18847887349020628, 'thresh': 2.0}, {'ineq': 'lt', 'dim': 12, 'alpha': 0.1522736899747682, 'thresh': 1.2}, {'ineq': 'gt', 'dim': 7, 'alpha': 0.15510870821690512, 'thresh': 1.2}, {'ineq': 'lt', 'dim': 5, 'alpha': 0.13536197353359397, 'thresh': 0.0}, {'ineq': 'lt', 'dim': 4, 'alpha': 0.12521587326132094, 'thresh': 28.799999999999997}, {'ineq': 'gt', 'dim': 11, 'alpha': 0.1334764812820768, 'thresh': 2.0}, {'ineq': 'lt', 'dim': 9, 'alpha': 0.14182243253771054, 'thresh': 4.0}, {'ineq': 'gt', 'dim': 14, 'alpha': 0.10264268449708046, 'thresh': 0.0}, {'ineq': 'lt', 'dim': 0, 'alpha': 0.11883732872109484, 'thresh': 1.0}, {'ineq': 'gt', 'dim': 4, 'alpha': 0.09879216527106671, 'thresh': 19.199999999999999}, {'ineq': 'lt', 'dim': 2, 'alpha': 0.12029960885056885, 'thresh': 36.719999999999999}, {'ineq': 'lt', 'dim': 3, 'alpha': 0.10846927663989193, 'thresh': 92.0}, {'ineq': 'lt', 'dim': 15, 'alpha': 0.09652967982091411, 'thresh': 0.0}, {'ineq': 'gt', 'dim': 3, 'alpha': 0.08958515309272004, 'thresh': 73.599999999999994}, {'ineq': 'lt', 'dim': 18, 'alpha': 0.09210361961272426, 'thresh': 8.900000000000004}, {'ineq': 'gt', 'dim': 16, 'alpha': 0.1046414221707963, 'thresh': 4.0}, {'ineq': 'lt', 'dim': 11, 'alpha': 0.09575457291711587, 'thresh': 3.2000000000000002}, {'ineq': 'gt', 'dim': 20, 'alpha': 0.09624217440331505, 'thresh': 0.0}, {'ineq': 'lt', 'dim': 17, 'alpha': 0.07859662885189679, 'thresh': 37.5}, {'ineq': 'lt', 'dim': 9, 'alpha': 0.071428636345507, 'thresh': 2.0}, {'ineq': 'gt', 'dim': 5, 'alpha': 0.07830753154662214, 'thresh': 2.0}, {'ineq': 'lt', 'dim': 4, 'alpha': 0.07606159074712804, 'thresh': 28.799999999999997}, {'ineq': 'gt', 'dim': 4, 'alpha': 0.08306752811081937, 'thresh': 19.199999999999999}, {'ineq': 'gt', 'dim': 7, 'alpha': 0.0830416741141177, 'thresh': 4.2000000000000002}, {'ineq': 'lt', 'dim': 3, 'alpha': 0.08893356802801233, 'thresh': 92.0}, {'ineq': 'gt', 'dim': 14, 'alpha': 0.07000509315417908, 'thresh': 3.0}, {'ineq': 'lt', 'dim': 7, 'alpha': 0.0769758358565884, 'thresh': 5.3000000000000005}, {'ineq': 'lt', 'dim': 18, 'alpha': 0.08507457442866707, 'thresh': 0.0}

```

这里输出了AdaBoost算法训练好的分类器的组合，我们只迭代了40次，也就是训练了40个弱分类器。最终，训练集的错误率为19.732%，测试集的错误率为19.403%，可以看到相对于Sklern的罗辑回归方法，错误率降低了很多。这个仅仅是我们训练40个弱分类器的结果，如果训练更多弱分类器，效果会更好。当分类器数量过多的时候，你会发现训练集错误率降低很多，但是测试集错误率提升了很多，这种现象就是**过拟合(overfitting)**。分类器对训练集的拟合效果好了普适性，只对训练集的分类效果好，这是我们不希望看到的。

2、使用Sklearn的AdaBoost

官方英文文档手册：[数据集下载](#)

sklearn.ensemble模块提供了很多集成方法，AdaBoost、Bagging、随机森林等。本文使用的是AdaBoostClassifier。

sklearn.ensemble: Ensemble Methods ¶

The `sklearn.ensemble` module includes ensemble-based methods for classification, regression and anomaly detection.

User guide: See the [Ensemble methods](#) section for further details.

<code>ensemble.AdaBoostClassifier</code> ([...])	An AdaBoost classifier.
<code>ensemble.AdaBoostRegressor</code> ([base_estimator, ...])	An AdaBoost regressor.
<code>ensemble.BaggingClassifier</code> ([base_estimator, ...])	A Bagging classifier.
<code>ensemble.BaggingRegressor</code> ([base_estimator, ...])	A Bagging regressor.
<code>ensemble.ExtraTreesClassifier</code> ([...])	An extra-trees classifier.
<code>ensemble.ExtraTreesRegressor</code> ([n_estimators, ...])	An extra-trees regressor.
<code>ensemble.GradientBoostingClassifier</code> ([loss, ...])	Gradient Boosting for classification.
<code>ensemble.GradientBoostingRegressor</code> ([loss, ...])	Gradient Boosting for regression.
<code>ensemble.IsolationForest</code> ([n_estimators, ...])	Isolation Forest Algorithm
<code>ensemble.RandomForestClassifier</code> ([...])	A random forest classifier.
<code>ensemble.RandomForestRegressor</code> ([...])	A random forest regressor.
<code>ensemble.RandomTreesEmbedding</code> ([...])	An ensemble of totally random trees.
<code>ensemble.VotingClassifier</code> (estimators[, ...])	Soft Voting/Majority Rule classifier for unfitted estimators.

让我们先看下AdaBoostClassifier这个函数，一共有5个参数：

sklearn.ensemble.AdaBoostClassifier

```
class sklearn.ensemble. AdaBoostClassifier (base_estimator=None, n_estimators=50, learning_rate=1.0,
algorithm='SAMME.R', random_state=None)
```

[\[source\]](#)

参数说明如下：

- **base_estimator**：可选参数，默认为DecisionTreeClassifier。理论上可以选择任何一个分类或者回归学习器，不过需要支持样本权重。我们常用CART决策树或者神经网络MLP。默认是决策树，即AdaBoostClassifier默认使用CART分类树DecisionTreeClassifier，而AdaBoostRegressor默认使用CART回归树DecisionTreeRegressor。另外有一个要注意的点是，如果我们选择的AdaBoostClassifier算法是SAMME.R，则我们的弱分类学习器支持概率预测，也就是在scikit-learn中弱分类学习器对应的预测方法除了predict还需要有predict_proba。
- **algorithm**：可选参数，默认为SAMME.R。scikit-learn实现了两种Adaboost分类算法，SAMME和SAMME.R。两者的主要区别是弱学习器权重。SAMME使用对样本集分类效果作为弱学习器权重，而SAMME.R使用了对样本集分类的预测概率大小来作为弱学习器权重。由于SAMME.R使用了连续值，迭代一般比SAMME快，因此AdaBoostClassifier的默认算法algorithm的值也是SAMME.R。我们一般使用默认的SAMME.R就够了，但如果是使用了SAMME.R，则弱分类学习器参数base_estimator必须限制使用支持概率预测的分类器。SAMME算法则没有这个限制。
- **n_estimators**：整数型，可选参数，默认为50。弱学习器的最大迭代次数，或者说最大的弱学习器的个数。一般来说n_estimators太小，容易欠拟合，n_estimators太大，又容易过拟合，一般选择一个适中的数值。默认是50。在实际调参的过程中，我们常常将n_estimators和下面介绍的参数learning_rate一起考虑。
- **learning_rate**：浮点型，可选参数，默认为1.0。每个弱学习器的权重缩减系数，取值范围为0到1，对于同样的训练集拟合效果，较小的v意味着更多的弱学习器的迭代次数。通常我们用步长和迭代最大次数一起来决定算法的拟合效果。所以这两个参数n_estimators和learning_rate要一起调参，可以从一个小一点的v开始调参，默认是1。
- **random_state**：整数型，可选参数，默认为None。如果RandomState的实例，random_state是随机数生成器；如果None，则随机数生成器是伪随机数生成器，使用np.random使用的RandomState实例。

了解这些，我们就可以开始编写代码了。完成上述代码相似的功能：

```
1 # -*-coding:utf-8 -*-
2 import numpy as np
3 from sklearn.ensemble import AdaBoostClassifier
4 from sklearn.tree import DecisionTreeClassifier
5 """
6 Author:
7     Jack Cui
8 Blog:
9     http://blog.csdn.net/c406495762
10 Zhihu:
11     https://www.zhihu.com/people/Jack--Cui/
12 Modify:
13     2017-10-11
14 """
15 def loadDataSet(fileName):
16     numFeat = len((open(fileName).readline().split('\t')))
17     dataMat = []; labelMat = []
18     fr = open(fileName)
19     for line in fr.readlines():
20         lineArr = []
21         curLine = line.strip().split('\t')
22         for i in range(numFeat - 1):
23             lineArr.append(float(curLine[i]))
24         dataMat.append(lineArr)
25         labelMat.append(float(curLine[-1]))
26     return dataMat, labelMat
27 if __name__ == '__main__':
28     dataArr, classLabels = loadDataSet('horseColicTraining2.txt')
29     testArr, testLabelArr = loadDataSet('horseColicTest2.txt')
30     bdt = AdaBoostClassifier(DecisionTreeClassifier(max_depth = 2), algorithm = "SAMME", n_estimators = 10)
31     bdt.fit(dataArr, classLabels)
32     predictions = bdt.predict(dataArr)
33     errArr = np.mat(np.ones((len(dataArr), 1)))
34     print('训练集的错误率:%.3f%%' % float(errArr[predictions != classLabels].sum() / len(dataArr) * 100))
35     predictions = bdt.predict(testArr)
36     errArr = np.mat(np.ones((len(testArr), 1)))
37     print('测试集的错误率:%.3f%%' % float(errArr[predictions != testLabelArr].sum() / len(testArr) * 100))
```

运行结果如下所示：

```
30
31 if __name__ == '__main__':
32     dataArr, classLabels = loadDataSet('horseColicTraining2.txt')
33     testArr, testLabelArr = loadDataSet('horseColicTest2.txt')
34     bdt = AdaBoostClassifier(DecisionTreeClassifier(max_depth = 2), algorithm = "SAMME", n_estimators =
35     bdt.fit(dataArr, classLabels)
36     predictions = bdt.predict(dataArr)
37     errArr = np.mat(np.ones((len(dataArr), 1)))
38     print('训练集的错误率:%.3f%%' % float(errArr[predictions != classLabels].sum() / len(dataArr) * 100))
39     predictions = bdt.predict(testArr)
40     errArr = np.mat(np.ones((len(testArr), 1)))
41     print('测试集的错误率:%.3f%%' % float(errArr[predictions != testLabelArr].sum() / len(testArr) * 100))

训练集的错误率:16.054%
测试集的错误率:17.910%
[Finished in 0.8s]
```

我们使用DecisionTreeClassifier作为使用的弱分类器，使用AdaBoost算法训练分类器。可以看到训练集的错误率为16.054%，测试集的错误率为：改n_estimators参数，你会发现跟我们自己写的代码，更改迭代次数的效果是一样的。n_estimators参数过大，会导致过拟合。

七、分类器性能评价

在之前的笔记中，讲了很多分类器。我们都是假设所有类别的分类代价是一样的。比如在逻辑回归那篇文章中，我们构建了一个用于检测患疝病的马系统。在那里，我们构建了分类器，但是并没有对分类后的情形加以讨论。假如某人给我们牵来一匹马，他希望我们能预测这匹马能否生存。我们说马会死可能会对马实施安乐死，而不是通过给马喂药来延缓其不可避免的死亡过程。我们的预测也许是错误的，马本来是可以继续活着的。毕竟，我们的分类器只（accuracy）。如果我们预测错误，那么我们将错杀一个如此昂贵的动物，更不要说人对马还存在情感上的依恋了。

再比如，如何过滤垃圾邮件呢？如果收件箱中会出现某些垃圾邮件，但合法邮件永远不会扔进垃圾邮件夹中，人们是否会满意呢？显然，我们可以出现的垃圾邮件，但是绝不能忍受，合法邮件被误扔如垃圾邮件夹中，万一这是一封女神or男神的表白信，这岂不是因此错过了一段旷世姻缘？

很多时候，不同类别的分类代价并不相等，这就是非均衡分类问题。我们将会考察一种新的分类器性能度量方法，而不再是简单的通过错误率进行评价技术来对上述非均衡问题下不同分类器性能进行可视化处理。

1、分类器性能度量指标

在之前，我们都是基于错误率来衡量分类器任务的成功程度的。错误率指的是在所有测试样本中错分的样本比例。实际上，这样的度量错误掩盖了样事实。在机器学习中，有一个普遍适用的称为混淆矩阵（confusion matrix）的工具，它可以帮助人们更好地了解分类中的错误。有这样一个关于在房子的动物类型的预测，这个预测的三个类问题的混淆矩阵如下图所示：

		预测结果		
		狗	猫	鼠
真实结果	狗	24	2	5
	猫	2	27	0
	鼠	4	2	30

利用混淆矩阵就可以更好地理解分类中的错误了。如果矩阵中的飞对角元素均为0，就会得到一个完美的分类器。

接下来，我们考虑另外一个混淆矩阵，这次的矩阵只针对一个简单的二类问题。混淆矩阵如下图所示：

		预测结果	
		+1	-1
真实结果	+1	真正例（TP）	伪反例（FN）
	-1	伪正例（FP）	真反例（TN）

可以看到，在这个二分类问题中，如果对一个正例正确地判为正例，那么就可以认为产生了一个真正例（True Positive，TP，也称真阳）；如果对一个反例，则认为产生了一个真反例（True Negative，TN，也称真阴）；如果对一个正例错误地判为反例，那么就可以认为产生了一个伪反例（False FN，为称假阴）；如果对一个反例错误地判为正例，则认为产生了一个伪正例（False Positive，FP，也称假阳）。

在分类中，当某个类别的重要性高于其他类别时，我们就可以来利用上述定义来定义出多个比错误率更好的指标。从混淆矩阵中，可以衍生出各种评(来自wiki)所示：

		Predicted condition			
		Total population	Predicted Condition positive	Predicted Condition negative	Prevalence $= \frac{\Sigma \text{Condition positive}}{\Sigma \text{Total population}}$
True condition	condition positive	True positive	False Negative (Type II error)	True positive rate (TPR), Sensitivity, Recall $= \frac{\Sigma \text{True positive}}{\Sigma \text{Condition positive}}$	False negative rate (FN Miss rate $= \frac{\Sigma \text{False negative}}{\Sigma \text{Condition positive}}$
	condition negative	False Positive (Type I error)	True negative	False positive rate (FPR), Fall-out $= \frac{\Sigma \text{False positive}}{\Sigma \text{Condition negative}}$	True negative rate (TN Specificity (SPC) $= \frac{\Sigma \text{True negative}}{\Sigma \text{Condition negative}}$
		Positive predictive value			

各个指标的定义及含义如下：

(1) Accuracy

模型的精度，即模型预测正确的个数/样本的总个数

$$Accuracy = \frac{TP+TN}{TP+FN+FP+TN}$$

一般情况下，模型的精度越高，说明模型的效果越好。

(2) Positive predictive value (PPV , Precision)

正确率，阳性预测值，在模型预测为正类的样本中，真正的正样本所占的比例。

$$Precision = \frac{TP}{TP+FP}$$

一般情况下，正确率越高，说明模型的效果越好。

(3) False discovery rate (FDR)

伪发现率，也是错误发现率，表示在模型预测为正类的样本中，真正的负类的样本所占的比例。

$$FDR = \frac{FP}{TP+FP}$$

一般情况下，错误发现率越小，说明模型的效果越好。

(4) False omission rate (FOR)

错误遗漏率，表示在模型预测为负类的样本中，真正的正类所占的比例。即评价模型“遗漏”掉的正类的多少。

$$FOR = \frac{FN}{FN+TN}$$

(5) Negative predictive value (NPV)

阴性预测值，在模型预测为负类的样本中，真正为负类的样本所占的比例。

$$NPV = \frac{TN}{FN+TN}$$

一般情况下，NPV越高，说明的模型的效果越好。

(6) True positive rate (TPR , Recall)

召回率，真正类率，表示的是，模型预测为正类的样本的数量，占总的正类样本数量的比值。

$$Recall = \frac{TP}{TP+FN}$$

一般情况下，Recall越高，说明有更多的正类样本被模型预测正确，模型的效果越好。

(7) False positive rate (FPR) , Fall-out

假正率，表示的是模型预测为正类的样本中，占模型负类样本数量的比值。

$$Fall-out = \frac{FP}{FP+TN}$$

一般情况下，假正类率越低，说明模型的效果越好。

(8) False negative rate (FNR) , Miss rate

假负类率，缺失率，模型预测为负类的样本中，是正类的数量，占真正类样本的比值。

$$FNR = \frac{FN}{FN+TN}$$

缺失值越小，说明模型的效果越好。

公式有些多？有些眼花缭乱，没关系，慢慢理解就好了。我们可以很容易构造一个高正确率或高召回率的分类器，但是很难同时保证两者成立。如果以正例为召回率，那么召回率达到百分之百而此时正确率很低。**构建一个同时使正确率和召回率最大的分类器是具有挑战性的。**

除了上述的评价指标，另一个用于度量分类中的非均衡的工具是ROC曲线（ROC curve），ROC代表接收者操作特征（receiver operating characteristic），早在二战期间由电气工程师构建雷达系统时使用过。

先运行个程序，我们看下结果，再听我细细道来，编写代码如下：

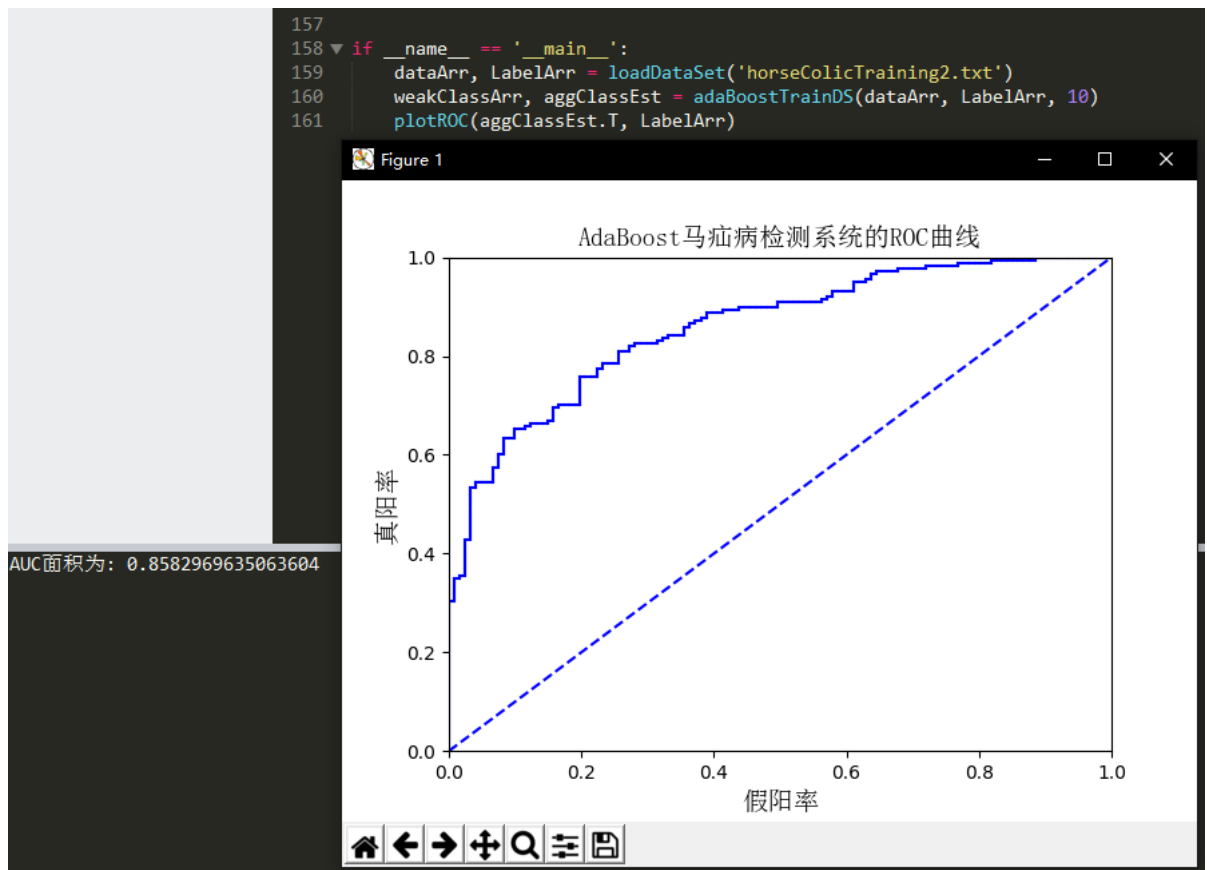
```
1 # -*-coding:utf-8 -*-
2 import numpy as np
3 import matplotlib.pyplot as plt
4 from matplotlib.font_manager import FontProperties
5
6 """
7 Author:
8     Jack Cui
9 Blog:
10     http://blog.csdn.net/c406495762
11 Zhihu:
12     https://www.zhihu.com/people/Jack--Cui/
13 Modify:
14     2017-10-10
15 """
16
17 def loadDataSet(fileName):
18     numFeat = len((open(fileName).readline().split('\t')))
19     dataMat = []; labelMat = []
20     fr = open(fileName)
21     for line in fr.readlines():
22         lineArr = []
23         curLine = line.strip().split('\t')
24         for i in range(numFeat - 1):
25             lineArr.append(float(curLine[i]))
26         dataMat.append(lineArr)
27         labelMat.append(float(curLine[-1]))
28
29     return dataMat, labelMat
30
31 def stumpClassify(dataMatrix, dimen, threshVal, threshIneq):
32     """
33     单层决策树分类函数
34     Parameters:
35         dataMatrix - 数据矩阵
36         dimen - 第dimen列，也就是第几个特征
37         threshVal - 阈值
38         threshIneq - 标志
39     Returns:
40         retArray - 分类结果
41     """
42     retArray = np.ones((np.shape(dataMatrix)[0],1)) #初始化retArray为1
43     if threshIneq == 'lt':
44         retArray[dataMatrix[:,dimen] <= threshVal] = -1.0 #如果小于阈值，则赋值为-1
45     else:
46         retArray[dataMatrix[:,dimen] > threshVal] = -1.0 #如果大于阈值，则赋值为-1
47     return retArray
48
49
50 def buildStump(dataArr, classLabels, D):
51     """
52     找到数据集上最佳的单层决策树
53     Parameters:
54         dataArr - 数据矩阵
55         classLabels - 数据标签
56         D - 样本权重
57     Returns:
58         bestStump - 最佳单层决策树信息
59         minError - 最小误差
60         bestClasEst - 最佳的分类结果
61     """
62     dataMatrix = np.mat(dataArr); labelMat = np.mat(classLabels).T
63     m, n = np.shape(dataMatrix)
64     numSteps = 10.0; bestStump = {}; bestClasEst = np.mat(np.zeros((m,1)))
65     minError = float('inf') #最小误差初始化为正无穷大
66     for i in range(n): #遍历所有特征
```

```

67     rangeMin = dataMatrix[:,i].min(); rangeMax = dataMatrix[:,i].max()           #找到特征中最小的值和最大值
68     stepSize = (rangeMax - rangeMin) / numSteps                                #计算步长
69     for j in range(-1, int(numSteps) + 1):
70         for unequal in ['lt', 'gt']:
71             threshVal = (rangeMin + float(j) * stepSize)                        #计算阈值
72             predictedVals = stumpClassify(dataMatrix, i, threshVal, unequal)    #计算分类结果
73             errArr = np.mat(np.ones((m,1)))                                    #初始化误差矩阵
74             errArr[predictedVals == labelMat] = 0                             #分类正确的,赋值为0
75             weightedError = D.T * errArr                                       #计算误差
76             # print("split: dim %d, thresh %.2f, thresh unequal: %s, the weighted error is %.3f" % (i, threshVal, unequal, weightedError))
77             if weightedError < minError:
78                 minError = weightedError
79                 bestClasEst = predictedVals.copy()
80                 bestStump['dim'] = i
81                 bestStump['thresh'] = threshVal
82                 bestStump['ineq'] = unequal
83     return bestStump, minError, bestClasEst
84
85 def adaBoostTrainDS(dataArr, classLabels, numIt = 40):
86     """
87     使用AdaBoost算法训练分类器
88     Parameters:
89         dataArr - 数据矩阵
90         classLabels - 数据标签
91         numIt - 最大迭代次数
92     Returns:
93         weakClassArr - 训练好的分类器
94         aggClassEst - 类别估计累计值
95     """
96     weakClassArr = []
97     m = np.shape(dataArr)[0]
98     D = np.mat(np.ones((m, 1)) / m)                                           #初始化权重
99     aggClassEst = np.mat(np.zeros((m,1)))
100    for i in range(numIt):
101        bestStump, error, classEst = buildStump(dataArr, classLabels, D)       #构建单层决策树
102        # print("D:",D.T)
103        alpha = float(0.5 * np.log((1.0 - error) / max(error, 1e-16)))         #计算弱学习算法权重alpha,使error不等于0,因为分母不能为0
104        bestStump['alpha'] = alpha
105        weakClassArr.append(bestStump)
106        # print("classEst: ", classEst.T)
107        expon = np.multiply(-1 * alpha * np.mat(classLabels).T, classEst)      #计算e的指数项
108        D = np.multiply(D, np.exp(expon))
109        D = D / D.sum()                                                         #根据样本权重公式,更新样本权重
110        #计算AdaBoost误差,当误差为0的时候,退出循环
111        aggClassEst += alpha * classEst
112        # print("aggClassEst: ", aggClassEst.T)
113        aggErrors = np.multiply(np.sign(aggClassEst) != np.mat(classLabels).T, np.ones((m,1))) #计算误差
114        errorRate = aggErrors.sum() / m
115        # print("total error: ", errorRate)
116        if errorRate == 0.0: break
117        #误差为0,退出循环
118    return weakClassArr, aggClassEst
119
120 def plotROC(predStrengths, classLabels):
121     """
122     绘制ROC
123     Parameters:
124         predStrengths - 分类器的预测强度
125         classLabels - 类别
126     Returns:
127         无
128     """
129     font = FontProperties(fname=r"c:\windows\fonts\simsun.ttc", size=14)
130     cur = (1.0, 1.0)
131     ySum = 0.0
132     numPosClas = np.sum(np.array(classLabels) == 1.0)
133     yStep = 1 / float(numPosClas)
134     xStep = 1 / float(len(classLabels) - numPosClas)
135
136     sortedIndicies = predStrengths.argsort()
137
138     fig = plt.figure()
139     fig.clf()
140     ax = plt.subplot(111)
141     for index in sortedIndicies.tolist()[0]:
142         if classLabels[index] == 1.0:
143             delX = 0; delY = yStep
144         else:
145             delX = xStep; delY = 0
146             ySum += cur[1]
147             ax.plot([cur[0], cur[0] - delX], [cur[1], cur[1] - delY], c = 'b')
148             cur = (cur[0] - delX, cur[1] - delY)
149             #高度累加
150             #绘制ROC
151             #更新绘制光标的位置
152     ax.plot([0, 1], [0, 1], 'b--')
153     plt.title('AdaBoost马痘病检测系统的ROC曲线', FontProperties = font)
154     plt.xlabel('假阳率', FontProperties = font)
155     plt.ylabel('真阳率', FontProperties = font)
156     ax.axis([0, 1, 0, 1])
157     print('AUC面积为:', ySum * xStep)
158     #计算AUC
159     plt.show()
160
161 if __name__ == '__main__':
162     dataArr, LabelArr = loadDataSet('horseColicTraining2.txt')
163     weakClassArr, aggClassEst = adaBoostTrainDS(dataArr, LabelArr, 10)
164     plotROC(aggClassEst.T, LabelArr)

```

运行结果如下所示：



我们可以看到有两个输出结果，一个是AUC面积，另一个ROC曲线图。

先解释ROC，图中的横坐标是伪正例的比例（假阳率=FP/（FP+TN）），而纵坐标是真正例的比例（真阳率=TP/（TP+FN））。ROC曲线给出的是假阳率和真阳率的变化情况。左下角的点所对应的将所有样例判为反例的情况，而右上角的点对应的则是将所有样例判为正例的情况。虚线给出的是随机线。

ROC曲线不但可以用于比较分类器，还可以基于成本效益（cost-versus-benefit）分析来做出决策。由于在不同的阈值下，不同的分类器的表现情况不同，因此以某种方式将它们组合起来或许更有意义。如果只是简单地观察分类器的错误率，那么我们就难以得到这种更深入的洞察效果了。

在理想的情况下，最佳的分类器应该尽可能地处于左上角，这就意味着分类器在假阳率很低的同时获得了很高的真阳率。例如在垃圾邮件的过滤中，识别所有的垃圾邮件，但没有将任何合法邮件误识别为垃圾邮件而放入垃圾邮件文件夹中。

对不同的ROC曲线进行比较的一个指标是曲线下的面积（Area Under the Curve，AUC）。AUC给出的是分类器的平均性能值，当然它并不能完全代替人的观察。一个完美分类器的AUC为1.0，而随机猜测的AUC则为0.5。

这个ROC曲线是怎么画的呢？

对于分类器而言，都有概率输出的功能，拿逻辑回归来举例，我们得到的是该样本属于正样本的概率和属于负样本的概率，属于正样本的概率大，那就判为正类，否则判为负类，那么实质上这里的阈值是0.5。

现在假设我们已经得到了所有样本的概率输出（属于正样本的概率），我们根据每个测试样本属于正样本的概率值从大到小排序。如下图所示：

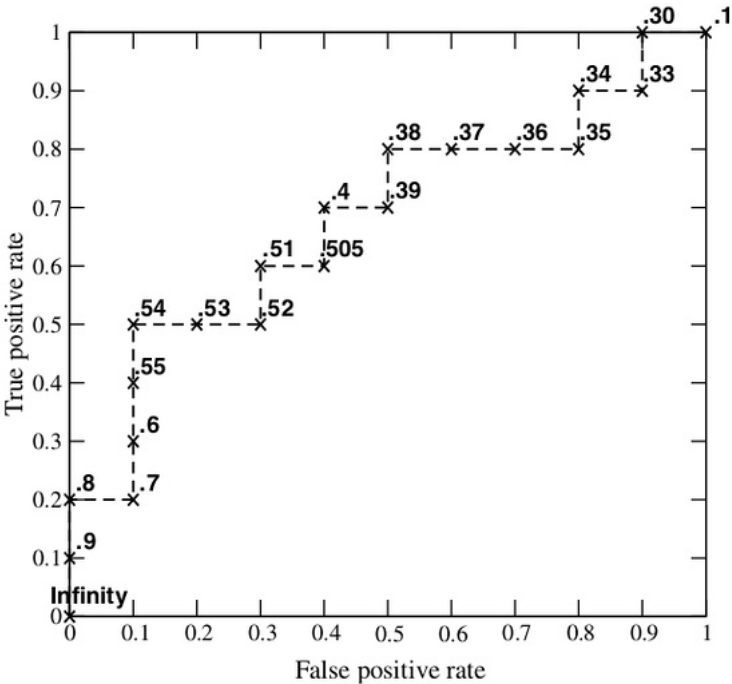
Inst#	Class	Score	Inst#	Class	Score
1	p	.9	11	p	.4
2	p	.8	12	n	.39
3	n	.7	13	p	.38
4	p	.6	14	n	.37
5	p	.55	15	n	.36
6	p	.54	16	n	.35
7	n	.53	17	p	.34
8	n	.52	18	n	.33
9	p	.51	19	p	.30
10	n	.505	20	n	.1

上图中共有20个测试样本，“Inst”一栏表示样本编号，“Class”一栏表示每个测试样本真正的标签（p表示正样本，n表示负样本），“Score”表示属于正样本的概率。其中，一共10个正样本，10个负样本。接下来，我们从高到低，依次将“Score”值作为阈值，当测试样本属于正样本的概率大于threshold时，我们认为它为正样本，否则为负样本。

举例来说：

- 当阈值为0.9时，样本1的“Score”值为0.9，大于等于阈值，那么样本1为正类，其余为负类，则 $TPR=TP/(TP+FN)=1/10=0.1$ ， $FPR=FP/(TP+FP)=0/10=0$ ；
- 当阈值为0.8时，样本1的“Score”值为0.9，样本2的“Score”值为0.8，大于等于阈值，那么样本1，2为正类，其余为负类，则 $TPR=2/10=0.2$ ， $FPR=0/10=0$ ；
- 当阈值为0.7时，样本1，2，3的“Score”值大于等于阈值，那么样本1，2，3为正类，其余为负类，则 $TPR=3/10=0.3$ ， $FPR=1/10=0.1$ ；
- 当阈值为0.6时，样本1，2，3，4的“Score”值大于等于阈值，那么样本1，2，3，4为正类，其余为负类，则 $TPR=4/10=0.4$ ， $FPR=1/10=0.1$ ；
- -----以此类推，此处省略一系列求解-----
- 当阈值为0.1时，所有样本的“Score”值大于等于阈值，那么所有样本均为正类，则 $TPR=10/10=1$ ， $FPR=10/10=1$ 。

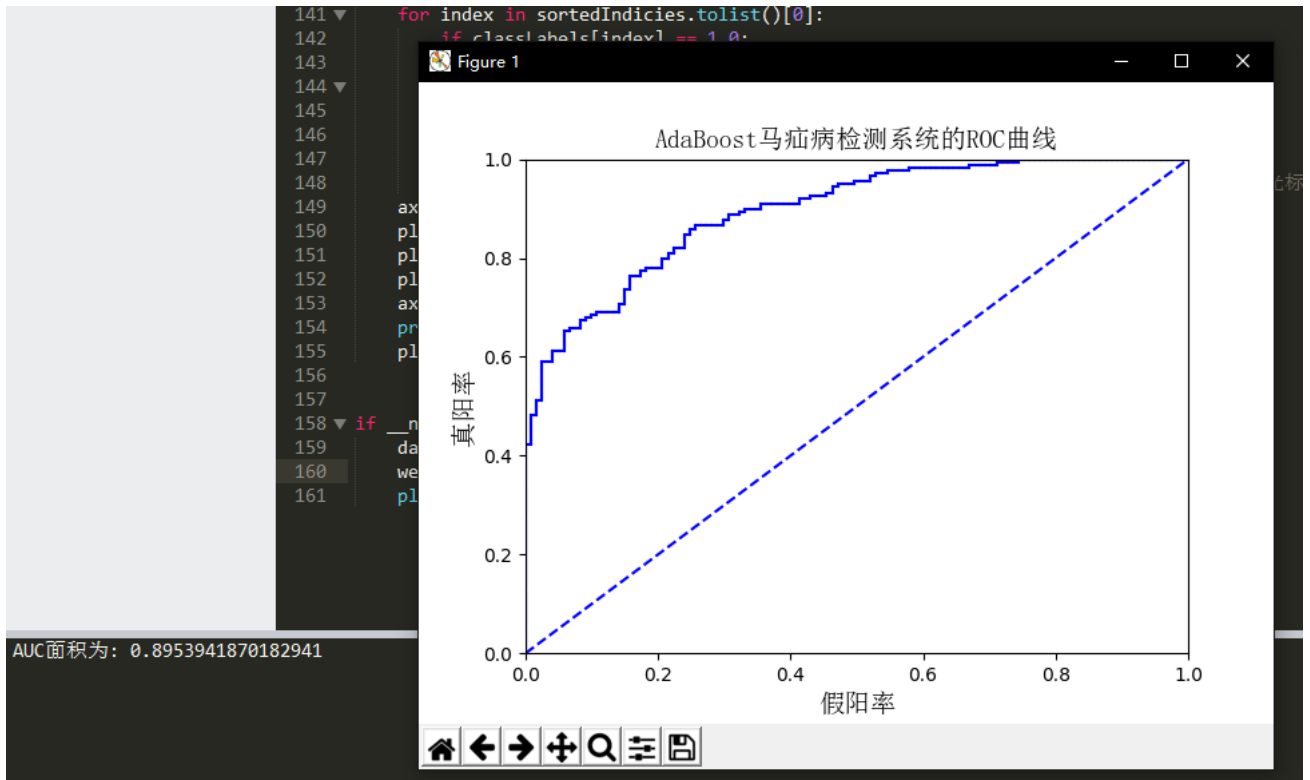
将它们的结果画出来，就构成了下图：



也许你已经发现了这样一个规律：多了一个TP，向Y轴移动一步，多了一个FP，向X轴移动一步。现在回头看一看绘制ROC曲线的程序吧，你会发现计算的，只不过区别在于，程序是从<1.0, 1.0>这个点开始画的。

AUC又是如何计算的呢？很简单，这些小矩形的宽度都是固定的xStep，因此先对所有矩形的高度进行累加，即ySum，最后面积就是ySum*xStep。

上面的ROC曲线绘制结果是在10个弱分类器下，AdaBoost算法性能的结果。我们将迭代次数改为50，也就是训练50个弱分类器，看下ROC曲线和A



你会发现，ROC曲线往左上角更靠拢了，并且AUC值增加了。也就表明，分类器效果更佳。

这就是ROC和AUC对与分类器评价指标，ROC越靠拢于左上角，分类器性能越好。同理AUC越接近于1，分类器性能越好。

八、总结

AdaBoost的优缺点：

- 优点：泛化错误率低，易编码，可以应用在大部分分类器上，无参数调整。
- 缺点：对离群点敏感。

其他：

- 集成方法通过组合多个分类器的分类结果，获得了比简单的单分类器更好的分类结果。本文只介绍了利用同一分类器的集成方法。除此之外，还有不同分类器的集成方法。
- 多个分类器组合可能会进一步凸显出单分类器的不足，比如过拟合问题。
- 本文主要介绍了boosting方法中最流行的一个称为AdaBoost的算法。
- 最后，本文介绍了一些分类器性能评价指标，召回率、ROC曲线、AUC等。
- 下篇文章开始讲解回归方法，欢迎届时捧场。
- 如有问题，请留言。如有错误，还望指正，谢谢！

PS：如果觉得本篇本章对您有所帮助，欢迎关注、评论、赞！

本文出现的所有代码和数据集，均可在我的github上下载，欢迎Follow、Star：[点击查看](#)

参考资料：

- [1] 简单易学的机器学习算法-AdaBoost：[点击查看](#)
- [2] Bagging和Boosting概念及区别：[点击查看](#)
- [3] 统计学习方法-CART, Bagging, Radom Forest, Boosting：[点击查看](#)
- [4] 十大经典数据挖掘算法AdaBoost：[点击查看](#)
- [5] 机器学习模型评价指标--混淆矩阵：[点击查看](#)
- [6] ROC和AUC介绍以及如何计算AUC：[点击查看](#)
- [7] scikit-learn AdaBoost类使用小结：[点击查看](#)



JackCui

关注人工智能及互联网的个人博客

[查看熊掌号](#)