

Boosting

Boosting的原理

Boosting方法使用同一组数据集进行反复学习，得到一系列简单模型。这些简单模型本身是弱学习器，且相互之间是强依赖的。然后将这些模型进行串行组合，得到一个预测性能更强大的机器学习模型。其主要是通过不断地减少偏差的形式来提高最终的预测结果。

Boosting方法最重要的两个问题是：

- 每一轮学习应该如何改变数据的概率分布
- 如何将各个弱分类器组合起来

Adaboost的原理

Adaboost的基本思想是：将弱学习器层层累加，在每一层训练的时候，对前一层基分类器分错的样本给予更高的权重（调整数据的分布）。在测试时，根据各层分类器的结果加权得到最终预测结果。

这个过程很类似于人类学习的过程，也即“**从错误中学习**”，我们学习新知识的过程往往是迭代的，第一遍学习的时候，我们会记住一部分知识，但是往往也会犯一些错误。在进一步学习时，我们会对之前犯错误的知识进行进一步的巩固，以减少此类错误的发生。如此不断循环往复，直到犯错误的次数减到很低的程度。

假设给定一个二分类的训练数据集： $T = \{(x_1, y_1), (x_2, y_2), \dots, (x_N, y_N)\}$ ，其中每个样本点由特征与类别组成。特征 $x_i \in \mathcal{X} \subseteq \mathbf{R}^n$ ，类别 $y_i \in \mathcal{Y} = \{-1, +1\}$ ， \mathcal{X} 是特征空间， \mathcal{Y} 是类别集合，输出最终分类器 $G(x)$ 。Adaboost算法如下：

(1) 初始化训练数据的分布： $D_1 = (w_{11}, \dots, w_{1i}, \dots, w_{1N})$ ， $w_{1i} = \frac{1}{N}$ ， $i = 1, 2, \dots, N$

(2) 对于 $m = 1, 2, \dots, M$

- 使用具有权值分布 D_m 的训练数据集进行学习，得到基本分类器： $G_m(x) : \mathcal{X} \rightarrow \{-1, +1\}$
- 计算 $G_m(x)$ 在训练集上的分类误差率
$$e_m = \sum_{i=1}^N P(G_m(x_i) \neq y_i) = \sum_{i=1}^N w_{mi} I(G_m(x_i) \neq y_i)$$
- 计算 $G_m(x)$ 的系数 $\alpha_m = \frac{1}{2} \log \frac{1-e_m}{e_m}$ ，这里的log是自然对数ln
- 更新训练数据集的权重分布

$$D_{m+1} = (w_{m+1,1}, \dots, w_{m+1,i}, \dots, w_{m+1,N})$$
$$w_{m+1,i} = \frac{w_{m,i}}{Z_m} \exp(-\alpha_m y_i G_m(x_i)), \quad i = 1, 2, \dots, N$$

这里的 Z_m 是规范化因子，使得 D_{m+1} 称为概率分布， $Z_m = \sum_{i=1}^N w_{mi} \exp(-\alpha_m y_i G_m(x_i))$

(3) 构建基本分类器的线性组合 $f(x) = \sum_{m=1}^M \alpha_m G_m(x)$ ，得到最终的分类器

$$G(x) = \text{sign}(f(x))$$
$$= \text{sign}\left(\sum_{m=1}^M \alpha_m G_m(x)\right)$$

Adaboost的实践

导入相关包

```

1 import pandas as pd
2 import numpy as np
3 from sklearn import datasets
4 from sklearn.model_selection import cross_val_score
5 from sklearn.tree import DecisionTreeClassifier
6 from sklearn.ensemble import AdaBoostClassifier
7 import matplotlib.pyplot as plt
8 plt.style.use("ggplot")

```

读入数据

```

1 iris = datasets.load_iris()
2 x = iris.data
3 y = iris.target
4 feature = iris.feature_names
5 data = pd.DataFrame(x, columns=feature)
6 data['target'] = y

```

构建模型

```

1 dt = DecisionTreeClassifier(criterion='entropy', random_state=1, max_depth=1)
2 ada =
  AdaBoostClassifier(base_estimator=dt, n_estimators=500, learning_rate=0.1, random_state=1)
3 models = [("dt", dt), ("adaboost", ada)]

```

为了方便对比，我们使用一层的决策树模型，该模型是一个较弱的分类器。

结果对比

```

1 def evaluate_model(model, x, y):
2     score = cross_val_score(model, x, y, scoring='accuracy', cv=5,
3     error_score='raise')
4     return score
5
6 for (name, model) in models:
7     score = evaluate_model(model, x, y)
8     print(f"Model:{name}; Mean: {score.mean():.3f}; Std: {score.std():.3f}")

```

Model:dt; Mean: 0.667; Std: 0.000

Model:adaboost; Mean: 0.947; Std: 0.034

可以看到虽然单层决策树的训练准确率较低，但是Adaboost算法能够有很高的性能提升。

可视化

为了更好地可视化，使用前两个维度的特征。

```

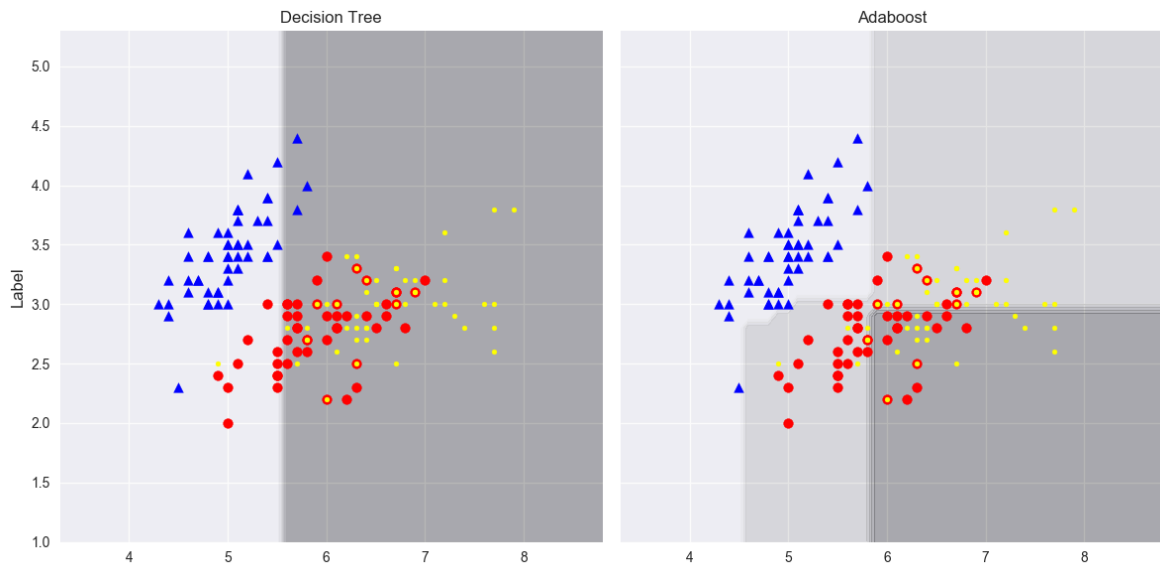
1 x = x[:, :2]
2 x_min = x[:, 0].min() - 1
3 x_max = x[:, 0].max() + 1
4 y_min = x[:, 1].min() - 1
5 y_max = x[:, 1].max() + 1
6 xx, yy = np.meshgrid(np.arange(x_min, x_max, 0.1), np.arange(y_min, y_max,
7 0.1))

```

```

7 f, axarr = plt.subplots(nrows=1, ncols=2, sharex='col', sharey='row', figsize=
  (12, 6))
8 for idx, clf, tt in zip([0, 1], [dt, ada], ['Decision Tree', 'Adaboost']):
9     clf.fit(X, y)
10    Z = clf.predict(np.c_[xx.ravel(), yy.ravel()])
11    Z = Z.reshape(xx.shape)
12    axarr[idx].contourf(xx, yy, Z, alpha=0.3)
13    axarr[idx].scatter(X[y == 0, 0], X[y == 0, 1], c='blue', marker='^')
14    axarr[idx].scatter(X[y == 1, 0], X[y == 1, 1], c='red', marker='o')
15    axarr[idx].scatter(X[y == 2, 0], X[y == 2, 1], c='yellow', marker='.')
16    axarr[idx].set_title(tt)
17 axarr[0].set_ylabel('Label', fontsize=12)
18 plt.tight_layout()
19 plt.show()

```



Adaboost模型的决策边界比单层决策树的决策边界要复杂的多。也就是说，Adaboost试图用增加模型复杂度而降低偏差的方式去减少总误差，但是过程中引入了方差，可能出现过拟合，因此在训练集和测试集之间的性能存在较大的差距。