

XGBoost 实验

PB18071477 敖旭扬

原理

给定 n 个样本、 m 个特征的二分类问题数据集 $D = \{(X_i, y_i)\} (|D| = n, X_i \in \mathbb{R}^m, y_i \in \{0, 1\})$ ，**XGBoost** 模型采用 K 次迭代的结果作为输出结果。 X_i 的输出为 \hat{y}_i ，最终预测结果为 p_i ，有

$$\hat{y}_i = \phi(X_i) = \sum_{k=1}^K f_k(X_i), \quad f_k \in F \quad (1)$$

$$p_i = \text{round}\left(\frac{1}{1 + e^{-\hat{y}_i}}\right) \quad (2)$$

其中 $F = \{f(X) = w_{q(X)}\} (q: \mathbb{R}^m \rightarrow \{1, 2, \dots, T\}, w \in \mathbb{R}^T)$ 表示 **XGBoost** 树结构空间集，各个变量的含义如下：

- q 表示树的结构，可以将样本映射到相应的叶子节点；
- T 表示 **XGBoost** 树叶子节点的数量；
- 每个 f_k 都对应独立的树结构 q 和权重 w 。

最终损失函数化简为

$$Obj^{(t)} = \sum_{j=1}^T [G_j w_j + \frac{1}{2} (H_j + \lambda) w_j^2] + \gamma T \quad (3)$$

其中

$$g_i = \frac{\partial l(y_i, \hat{y}_i^{(t-1)})}{\partial \hat{y}_i^{(t-1)}}, \quad h_i = \frac{\partial^2 l(y_i, \hat{y}_i^{(t-1)})}{\partial^2 \hat{y}_i^{(t-1)}} \quad (4)$$

$$I_j = \{i | q(X_i) = j\}$$

$$G_j = \sum_{i \in I_j} g_i, \quad H_j = \sum_{i \in I_j} h_i \quad (5)$$

$Obj^{(t)}$ 对 w_j 求导等于 0 可得

$$w_j^* = -\frac{G_j}{H_j + \lambda} \quad (6)$$

所以接下来需要不断寻找划分方式来构造出树的结构，可以对于数据集，遍历所有划分属性和分界点，求出

$$Gain = \frac{1}{2} \left[\frac{G_L^2}{H_L + \lambda} + \frac{G_R^2}{H_R + \lambda} - \frac{(G_L + G_R)^2}{H_L + H_R + \lambda} \right] \quad (7)$$

按照取得 $Gain$ 最大值的方式进行划分，再对划分后的两个子集继续按照这样的规则划分，同时树的层次增加 1。当树的层次达到最大值或当前集合小于最小划分块大小时停止划分得到一个叶子结点，按照 (6) 式计算出该叶子的权值。训练完当前树，更新 \hat{y}_i 后再训练下一棵树，最终得到包含 K 棵树的 **XGBoost** 模型，据此可用该模型对数据集进行预测，输出训练集和测试集的精度。

在本实现中，损失函数使用的是 *logistic* 损失，即


```

25         h_l = np.sum(XypGH[:split_i, -1])
26         h_r = np.sum(XypGH[split_i:, -1])
27         gain = 0.5*(np.square(g_l)/(h_l+self.la)
28                 + np.square(g_r)/(h_r+self.la)
29                 - np.square(g_l+g_r)/(h_l+h_r+self.la))
30         - self.gamma
31         if gain >= max_gain and gain > self.min_gain:
32             max_gain = gain
33             feature = f
34             threshold = fv
35             split_point = split_i
36         if max_gain > self.min_gain:
37             XypGH = XypGH[np.lexsort(XypGH.T[feature, None])]
38             LXypGH = XypGH[:split_point]
39             RXypGH = XypGH[split_point:]
40             left = self.build_tree(LXypGH, cur_depth + 1)
41             right = self.build_tree(RXypGH, cur_depth + 1)
42             return self.Node(feature, threshold, left, right)
43     g = np.sum(XypGH[:, -2])
44     h = np.sum(XypGH[:, -1])
45     leaf_value = self.eta*(-g / (h + self.la))
46     return self.Node(value=leaf_value)
47
48     def predict_value(self, x, tree):
49         if tree.value is not None:
50             return tree.value
51         feature_value = x[tree.feature]
52         if feature_value < tree.threshold:
53             return self.predict_value(x, tree.left)
54         else:
55             return self.predict_value(x, tree.right)
56
57     def predict(self, X):
58         y_pred = []
59         for x in X:
60             y_pred.append(self.predict_value(x, self.root))
61         p = np.array(y_pred)
62         return p
63

```

完整实验源码见压缩包中的[XGBoost.py](#)。

运算结果

实例

在主函数中使用默认超参数(详见源码)调用下面的实例

```

1 # 训练模型
2 xgbt = XGBoost()
3 xgbt.fit(x_train, y_train)
4
5 # 训练集和测试集精度
6 p_train = xgbt.predict(x_train)
7 p_test = xgbt.predict(x_test)
8 print("训练集精度为: {:.2f} %".format(
9     (p_train == y_train).mean() * 100))
10 print("测试集精度为: {:.2f} %".format(
11     (p_test == y_test).mean() * 100))

```

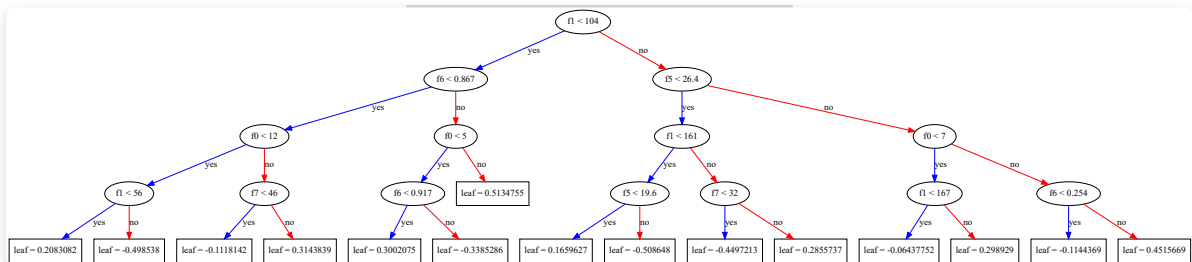
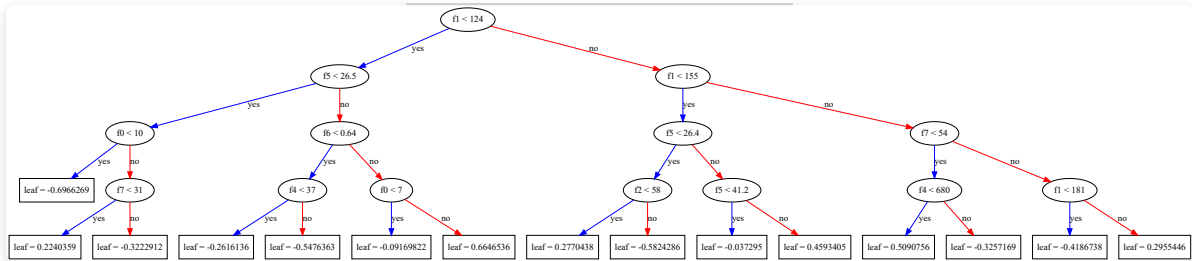
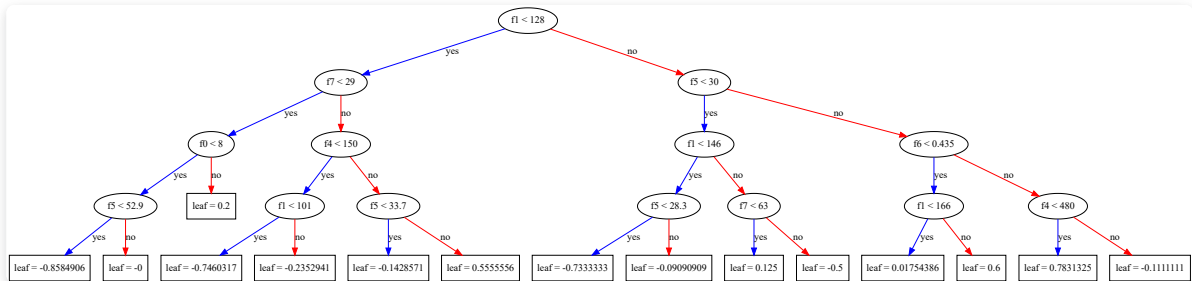
命令行输出结果为

```

1 训练集精度为: 83.58 %
2 测试集精度为: 76.47 %

```

训练结果



预测效果

在测试集上的预测精度为 76.47%。

总结

题目要求的 Baseline 为

- 1 | 测评指标：精度值，正确预测占整体的比例
- 2 | 测试集精度：0.7

我训练出的**XGBoost模型**测试集精度为 76.47%，性能达标。