

# 作业 3：决策树与提升算法

## 2 决策树（40pt）

### 2.1 ID3 决策树算法（15pt）

1. 虽然表中没有列出数据集中的具体实例，但已经足以构建一棵 ID3 决策树的根节点（决策树桩）。请通过计算说明应该如何构建。

计算按是否晴天分类的信息增益

$$H(\mathcal{D}) = -\frac{35}{50} \log \frac{35}{50} - \frac{15}{50} \log \frac{15}{50} \approx 0.8813$$

$$H(\mathcal{D}_1) = -\frac{28}{34} \log \frac{28}{34} - \frac{6}{34} \log \frac{6}{34} \approx 0.6723$$

$$H(\mathcal{D}_2) = -\frac{7}{16} \log \frac{7}{16} - \frac{9}{16} \log \frac{9}{16} \approx 0.9887$$

$$H(\mathcal{D}) - \frac{|\mathcal{D}_1|}{|\mathcal{D}|} H(\mathcal{D}_1) - \frac{|\mathcal{D}_2|}{|\mathcal{D}|} H(\mathcal{D}_2) = 0.8813 - \frac{34}{50} \times 0.6723 - \frac{16}{50} \times 0.9887 \approx 0.1077$$

计算按是否有雪分类的信息增益

$$H(\mathcal{D}) = -\frac{35}{50} \log \frac{35}{50} - \frac{15}{50} \log \frac{15}{50} \approx 0.8813$$

$$H(\mathcal{D}_1) = -\frac{16}{18} \log \frac{16}{18} - \frac{2}{18} \log \frac{2}{18} \approx 0.5033$$

$$H(\mathcal{D}_2) = -\frac{19}{32} \log \frac{19}{32} - \frac{13}{32} \log \frac{13}{32} \approx 0.9745$$

$$H(\mathcal{D}) - \frac{|\mathcal{D}_1|}{|\mathcal{D}|} H(\mathcal{D}_1) - \frac{|\mathcal{D}_2|}{|\mathcal{D}|} H(\mathcal{D}_2) = 0.8813 - \frac{18}{50} \times 0.5033 - \frac{32}{50} \times 0.9745 \approx 0.0764$$

因此，根节点应该按是否晴天分类。

2. 使用最小错误替换 ID3 算法信息增益公式中的熵函数，构建一棵完整的决策树（展现过程）。

首先构造根节点。

$$\text{MinError}(\mathcal{D}) = \min \left\{ \frac{11}{24}, \frac{13}{24} \right\} = \frac{11}{24}$$

计算按颜色分类的信息增益

$$\text{MinError}(\mathcal{D}_1) = \min \left\{ \frac{5}{13}, \frac{8}{13} \right\} = \frac{5}{13}$$

$$\text{MinError}(\mathcal{D}_2) = \min \left\{ \frac{6}{11}, \frac{5}{11} \right\} = \frac{5}{11}$$

$$\text{MinError}(\mathcal{D}) - \frac{|\mathcal{D}_1|}{|\mathcal{D}|} \text{MinError}(\mathcal{D}_1) - \frac{|\mathcal{D}_2|}{|\mathcal{D}|} \text{MinError}(\mathcal{D}_2) = \frac{11}{24} - \frac{13}{24} \times \frac{5}{13} - \frac{11}{24} \times \frac{5}{11} = \frac{1}{24}$$

计算按大小分类的信息增益

$$\text{MinError}(\mathcal{D}_1) = \min \left\{ \frac{7}{18}, \frac{11}{18} \right\} = \frac{7}{18}$$

$$\text{MinError}(\mathcal{D}_2) = \min \left\{ \frac{2}{6}, \frac{4}{6} \right\} = \frac{2}{6}$$

$$\text{MinError}(\mathcal{D}) - \frac{|\mathcal{D}_1|}{|\mathcal{D}|}\text{MinError}(\mathcal{D}_1) - \frac{|\mathcal{D}_2|}{|\mathcal{D}|}\text{MinError}(\mathcal{D}_2) = \frac{11}{24} - \frac{18}{24} \times \frac{7}{18} - \frac{6}{24} \times \frac{2}{6} = \frac{1}{12}$$

计算按动作分类的信息增益

$$\text{MinError}(\mathcal{D}_1) = \min \left\{ \frac{5}{12}, \frac{7}{12} \right\} = \frac{5}{12}$$

$$\text{MinError}(\mathcal{D}_2) = \min \left\{ \frac{6}{12}, \frac{6}{12} \right\} = \frac{6}{12}$$

$$\text{MinError}(\mathcal{D}) - \frac{|\mathcal{D}_1|}{|\mathcal{D}|}\text{MinError}(\mathcal{D}_1) - \frac{|\mathcal{D}_2|}{|\mathcal{D}|}\text{MinError}(\mathcal{D}_2) = \frac{11}{24} - \frac{12}{24} \times \frac{5}{12} - \frac{12}{24} \times \frac{6}{12} = 0$$

计算按年龄分类的信息增益

$$\text{MinError}(\mathcal{D}_1) = \min \left\{ \frac{6}{12}, \frac{6}{12} \right\} = \frac{6}{12}$$

$$\text{MinError}(\mathcal{D}_2) = \min \left\{ \frac{5}{12}, \frac{7}{12} \right\} = \frac{5}{12}$$

$$\text{MinError}(\mathcal{D}) - \frac{|\mathcal{D}_1|}{|\mathcal{D}|}\text{MinError}(\mathcal{D}_1) - \frac{|\mathcal{D}_2|}{|\mathcal{D}|}\text{MinError}(\mathcal{D}_2) = \frac{11}{24} - \frac{12}{24} \times \frac{6}{12} - \frac{12}{24} \times \frac{5}{12} = 0$$

因此，根节点应该按大小分类。

对于“小”的那部分，显然可以根据颜色或年龄分类。

对于“大”的那部分

$$\text{MinError}(\mathcal{D}) = \min \left\{ \frac{7}{18}, \frac{11}{18} \right\} = \frac{7}{18}$$

计算按颜色分类的信息增益

$$\text{MinError}(\mathcal{D}_1) = \min \left\{ \frac{5}{13}, \frac{8}{13} \right\} = \frac{5}{13}$$

$$\text{MinError}(\mathcal{D}_2) = \min \left\{ \frac{2}{5}, \frac{3}{5} \right\} = \frac{2}{5}$$

$$\text{MinError}(\mathcal{D}) - \frac{|\mathcal{D}_1|}{|\mathcal{D}|}\text{MinError}(\mathcal{D}_1) - \frac{|\mathcal{D}_2|}{|\mathcal{D}|}\text{MinError}(\mathcal{D}_2) = \frac{7}{18} - \frac{13}{18} \times \frac{5}{13} - \frac{5}{18} \times \frac{2}{5} = 0$$

计算按动作分类的信息增益

$$\text{MinError}(\mathcal{D}_1) = \min \left\{ \frac{5}{12}, \frac{7}{12} \right\} = \frac{5}{12}$$

$$\text{MinError}(\mathcal{D}_2) = \min \left\{ \frac{2}{6}, \frac{4}{6} \right\} = \frac{2}{6}$$

$$\text{MinError}(\mathcal{D}) - \frac{|\mathcal{D}_1|}{|\mathcal{D}|}\text{MinError}(\mathcal{D}_1) - \frac{|\mathcal{D}_2|}{|\mathcal{D}|}\text{MinError}(\mathcal{D}_2) = \frac{7}{18} - \frac{12}{18} \times \frac{5}{12} - \frac{6}{18} \times \frac{2}{6} = 0$$

计算按年龄分类的信息增益

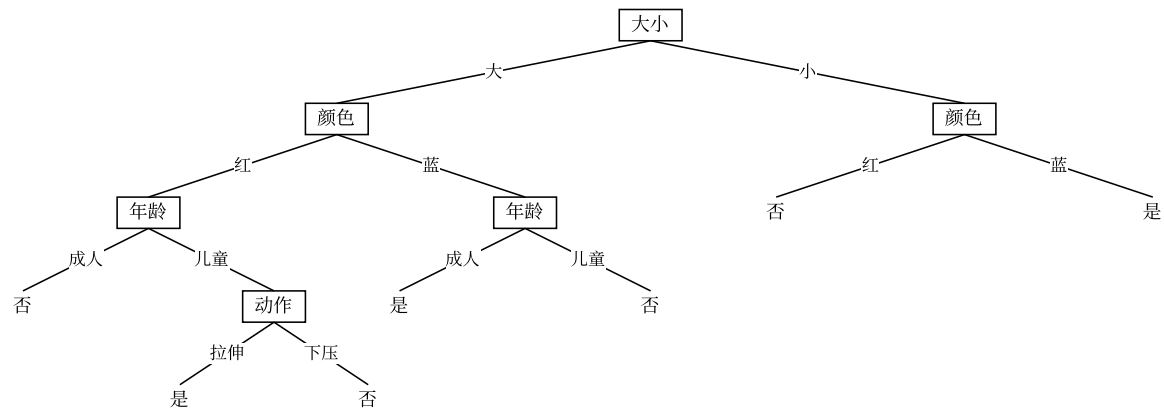
$$\text{MinError}(\mathcal{D}_1) = \min \left\{ \frac{2}{8}, \frac{6}{8} \right\} = \frac{2}{8}$$

$$\text{MinError}(\mathcal{D}_2) = \min \left\{ \frac{5}{10}, \frac{5}{10} \right\} = \frac{5}{10}$$

$$\text{MinError}(\mathcal{D}) - \frac{|\mathcal{D}_1|}{|\mathcal{D}|} \text{MinError}(\mathcal{D}_1) - \frac{|\mathcal{D}_2|}{|\mathcal{D}|} \text{MinError}(\mathcal{D}_2) = \frac{7}{18} - \frac{8}{18} \times \frac{2}{8} - \frac{10}{18} \times \frac{5}{10} = 0$$

由于三者信息增益相同，不妨选择颜色分类。

以此类推，可以构造出完整的决策树。



3. ID3 算法是否总能保证生成一颗“最优”的决策树？这里，“最优”的定义是：决策树能够最好的拟合训练数据，且深度最小。若能，请说明原因；若不能，请举出反例。

ID3 算法不能保证生成一颗全局最优的决策树，它是一种贪心算法，只是在局部选择最优属性分裂点，无法确保获得整体最优解。

我们可以考虑如下数据集

A	B	C	Label
0	0	0	1
0	0	1	0
0	1	0	0
0	1	1	1
1	0	0	1
1	0	1	0
1	1	0	0
1	1	1	1

在根节点处，由于选择 A、B、C 三个属性的信息增益相同，ID3 算法会随机选择一个属性进行分裂。假设选择 A 属性进行分裂，此时决策树深度为 3。事实上，如果选择 B 属性或 C 属性进行分裂，都可以得到深度为 2 的决策树。所以 ID3 算法不能保证生成一颗“最优”的决策树。

## 2.2 代码实验（25pt）

1. 根据熵的定义，完成 `tree.py` 中 `compute_entropy` 函数。

```
1 def compute_entropy(label_array):
2     """
3     Calculate the entropy of given label list
4
5     :param label_array: a numpy array of labels shape = (n, 1)
```

```

6         :return entropy: entropy value
7         """
8         # Your code goes here (~6 lines)
9         # TODO 2.3.1
10        _, label_counts = np.unique(label_array, return_counts=True)
11        prob = label_counts / len(label_array)
12        entropy = -np.sum(prob * np.log2(prob))
13        return entropy

```

2. 根据基尼系数的定义，完成 `tree.py` 中 `compute_gini` 函数。

```

1  def compute_gini(label_array):
2      """
3      Calculate the gini index of label list
4
5      :param label_array: a numpy array of labels shape = (n, 1)
6      :return gini: gini index value
7      """
8      # Your code goes here (~6 lines)
9      # TODO 2.3.2
10     _, label_counts = np.unique(label_array, return_counts=True)
11     prob = label_counts / len(label_array)
12     gini = 1 - np.sum(prob ** 2)
13     return gini

```

3. 补全 `tree.py` 中 `DecisionTree` 类的 `fit` 函数。提示：递归调用决策树的构造与 `fit` 函数。

```

1  def fit(self, X, y=None):
2      """
3      This should fit the tree classifier by setting the values self.is_leaf,
4      self.split_id (the index of the feature we want to split on, if we're
5      splitting),
6      self.split_value (the corresponding value of that feature where the
7      split is),
8      and self.leaf_value, which is the prediction value if the tree is a
9      leaf node. If we
10     are splitting the node, we should also init self.left and self.right to
11     be DecisionTree
12     objects corresponding to the left and right subtrees. These subtrees
13     should be fit on
14     the data that fall to the left and right, respectively, of
15     self.split_value.
16     This is a recursive tree building procedure.
17
18     :param X: a numpy array of training data, shape = (n, m)
19     :param y: a numpy array of labels, shape = (n, 1)
20
21     :return self
22     """
23     # If depth is max depth turn into leaf
24     if self.depth == self.max_depth:
25         self.is_leaf = True
26         self.leaf_value = self.leaf_value_estimator(y)
27         return self

```

```

23     # If reach minimum sample size turn into leaf
24     if len(y) ≤ self.min_sample:
25         self.is_leaf = True
26         self.leaf_value = self.leaf_value_estimator(y)
27         return self
28
29     # If not is_leaf, i.e in the node, we should create left and right
subtree
30     # But First we need to decide the self.split_id and self.split_value
that minimize loss
31     # Compare with constant prediction of all X
32     best_split_value = None
33     best_split_id = None
34     best_loss = self.split_loss_function(y)
35     best_left_X = None
36     best_right_X = None
37     best_left_y = None
38     best_right_y = None
39     # Concatenate y into X for sorting together
40     X = np.concatenate([X, y], 1)
41     for i in range(X.shape[1] - 1):
42         # Note: The last column of X is y now
43         X = np.array(sorted(X, key=lambda x: x[i]))
44         for split_pos in range(len(X) - 1):
45             # :split_pos+1 will include the split_pos data in left_X
46             left_X = X[:split_pos + 1, :-1]
47             right_X = X[split_pos + 1:, :-1]
48             # you need left_y to be in (n,1) i.e (-1,1) dimension
49             left_y = X[:split_pos + 1, -1].reshape(-1, 1)
50             right_y = X[split_pos + 1:, -1].reshape(-1, 1)
51             left_loss = len(left_y) * self.split_loss_function(left_y) /
len(y)
52             right_loss = len(right_y) * self.split_loss_function(right_y) /
len(y)
53             # If any choice of splitting feature and splitting position
results in better loss
54             # record following information and discard the old one
55             if ((left_loss + right_loss) < best_loss):
56                 best_split_value = X[split_pos, i]
57                 best_split_id = i
58                 best_loss = left_loss + right_loss
59                 best_left_X = left_X
60                 best_right_X = right_X
61                 best_left_y = left_y
62                 best_right_y = right_y
63
64     # Condition when you have a split position that results in better loss
65     # Your code goes here (~10 lines)
66     # TODO 2.3.3
67     if best_split_id ≠ None:
68         # build child trees and set splitting info
69         self.split_id = best_split_id
70         self.split_value = best_split_value
71         self.left = DecisionTree(self.split_loss_function,
self.leaf_value_estimator, self.depth+1, self.min_sample,
max_depth=self.max_depth)

```

```

72         self.right = DecisionTree(self.split_loss_function,
self.leaf_value_estimator, self.depth+1, self.min_sample,
max_depth=self.max_depth)
73         self.left.fit(best_left_X, best_left_y)
74         self.right.fit(best_right_X, best_right_y)
75     else:
76         # set value
77         self.is_leaf = True
78         self.leaf_value = self.leaf_value_estimator(y)
79
80     return self

```

4. 完成 `tree.py` 中 `mean_absolute_deviation_around_median` 函数。

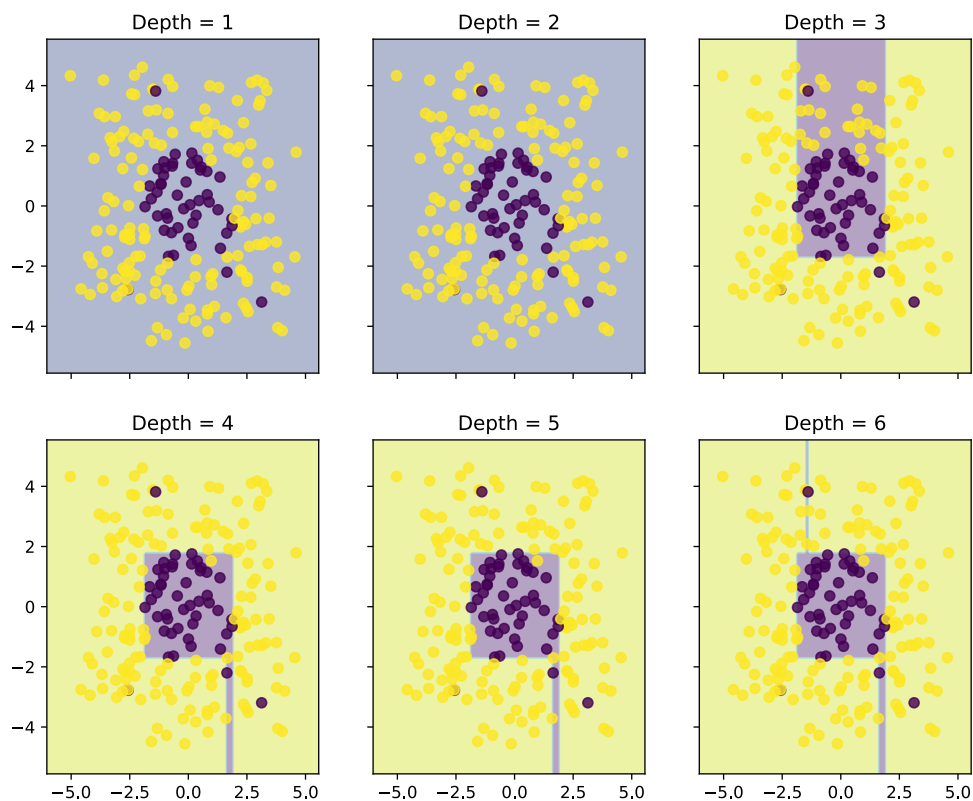
```

1  # Regression Tree Specific Code
2  def mean_absolute_deviation_around_median(y):
3      """
4          Calculate the mean absolute deviation around the median of a given
target list
5
6          :param y: a numpy array of targets shape = (n, 1)
7          :return mae
8          """
9          # Your code goes here (~3 lines)
10         # TODO 2.3.4
11         mae = np.mean(np.abs(y - np.median(y)))
12         return mae

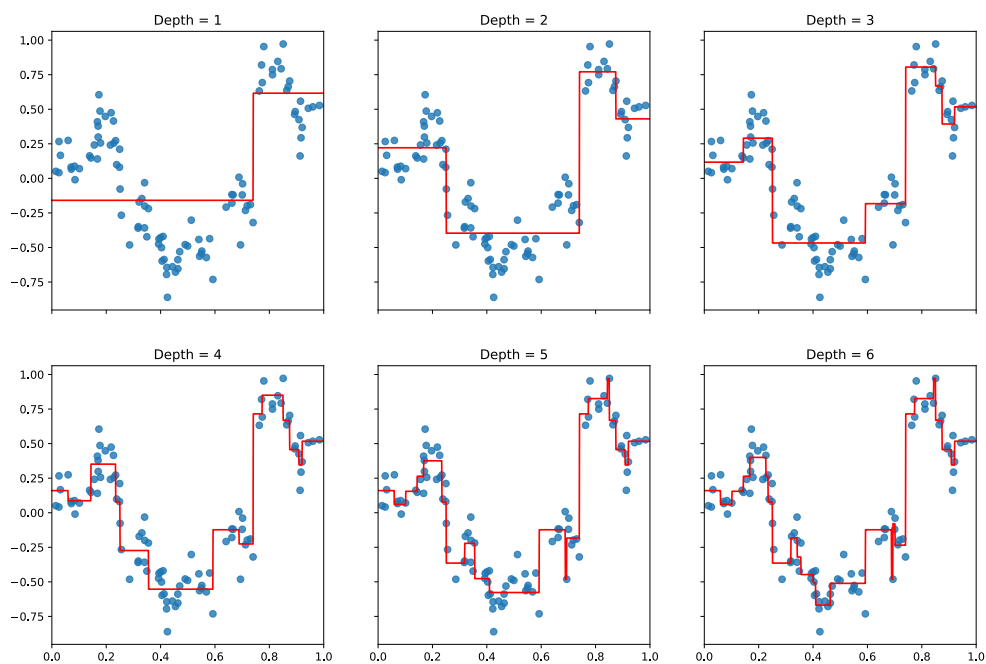
```

5. 运行 `tree.py`，在实验文档中记录决策树在不同数据集上运行的结果，包括 (a) `DT_entropy.pdf`，使用决策树在二分类问题上的结果。(b) `DT_regression.pdf`，使用决策树在回归问题上的结果。并简要描述实验现象（例如超参数对于决策树的影响）。

`DT_entropy.pdf` 的结果如下



DT\_regression.pdf 的结果如下



可以观察到，在分类任务中，决策树的深度较小时（Depth = 1 或 Depth = 2），决策树完全无法进行分类；随着深度的增加，决策树的分类效果逐渐提升，但是深度过大时（Depth = 6）会出现过拟合现象。在回归任务中，决策树的深度较小时（Depth = 1），决策树的拟合效果较差，无法反映数据的变化；随着深度的增加，决策树对数据的拟合效果逐渐提升，但是深度过大时（Depth = 5 或 Depth = 6）会出现过拟合现象，即有很多小毛刺。

### 3 提升算法 (60pt+10pt)

#### 3.1 弱分类器的更新保证 (10pt)

1. 通过计算证明:  $Z_t \doteq \sum_{i=1}^n D_t(i) \exp(-\alpha_t y_i h_t(\mathbf{x}_i)) = 2[\epsilon_t(1 - \epsilon_t)]^{\frac{1}{2}}$ 。

归一化因子  $Z_t$  的定义为

$$Z_t = \sum_{i=1}^n D_t(i) \exp(-\alpha_t y_i h_t(\mathbf{x}_i))$$

其中, 当分类正确, 即  $y_i = h_t(\mathbf{x}_i)$  时, 有  $y_i h_t(\mathbf{x}_i) = 1$ ; 当分类错误, 即  $y_i \neq h_t(\mathbf{x}_i)$  时, 有  $y_i h_t(\mathbf{x}_i) = -1$ 。故而可以将  $Z_t$  分解为两部分:

$$Z_t = \sum_{y_i=h_t(\mathbf{x}_i)} D_t(i) \exp(-\alpha_t) + \sum_{y_i \neq h_t(\mathbf{x}_i)} D_t(i) \exp(\alpha_t)$$

由于有

$$\epsilon_t = \Pr_{i \sim D_t} [h_t(\mathbf{x}_i) \neq y_i] = \sum_{y_i \neq h_t(\mathbf{x}_i)} D_t(i)$$

故而可以将  $Z_t$  进一步写成

$$Z_t = \epsilon_t \exp(\alpha_t) + (1 - \epsilon_t) \exp(-\alpha_t)$$

考虑到

$$\alpha_t \leftarrow \frac{1}{2} \log \left( \frac{1 - \epsilon_t}{\epsilon_t} \right)$$

故而有

$$Z_t = \epsilon_t \exp \left( \frac{1}{2} \log \left( \frac{1 - \epsilon_t}{\epsilon_t} \right) \right) + (1 - \epsilon_t) \exp \left( -\frac{1}{2} \log \left( \frac{1 - \epsilon_t}{\epsilon_t} \right) \right) = 2[\epsilon_t(1 - \epsilon_t)]^{\frac{1}{2}}$$

2. 证明:  $h_t$  关于分布  $D_{t+1}$  的错误率正好为  $\frac{1}{2}$ , 即对任意  $1 \leq t < T$ :

$\sum_{i=1}^n D_{t+1}(i) \mathbb{1}_{[y_i \neq h_t(\mathbf{x}_i)]} = \frac{1}{2}$ , 并据此说明, 对任意  $1 \leq t < T$ ,  $t+1$  步选取的弱分类器  $h_{t+1}$  不会与  $h_t$  相同。

显然

$$\sum_{i=1}^n D_{t+1}(i) \mathbb{1}_{[y_i \neq h_t(\mathbf{x}_i)]} + \sum_{i=1}^n D_{t+1}(i) \mathbb{1}_{[y_i = h_t(\mathbf{x}_i)]} = 1$$

考虑到

$$D_{t+1}(i) \leftarrow \frac{D_t(i) \exp(-\alpha_t y_i h_t(\mathbf{x}_i))}{Z_t}$$

其中, 当分类正确, 即  $y_i = h_t(\mathbf{x}_i)$  时, 有  $y_i h_t(\mathbf{x}_i) = 1$ ; 当分类错误, 即  $y_i \neq h_t(\mathbf{x}_i)$  时, 有  $y_i h_t(\mathbf{x}_i) = -1$ 。故而可以将上式改写成:

$$\sum_{y_i \neq h_t(\mathbf{x}_i)} \frac{D_t(i) \exp(\alpha_t)}{Z_t} + \sum_{y_i = h_t(\mathbf{x}_i)} \frac{D_t(i) \exp(-\alpha_t)}{Z_t} = 1$$



又因为

$$\epsilon_t = \Pr_{i \sim D_t} [h_t(\mathbf{x}_i) \neq y_i] = \sum_{y_i \neq h_t(\mathbf{x}_i)} D_t(i)$$

和

$$\alpha_t \leftarrow \frac{1}{2} \log \left( \frac{1 - \epsilon_t}{\epsilon_t} \right)$$

以及上一题证明的

$$Z_t = 2[\epsilon_t(1 - \epsilon_t)]^{\frac{1}{2}}$$

故而有

$$\sum_{y_i \neq h_t(\mathbf{x}_i)} \frac{D_t(i) \exp(\alpha_t)}{Z_t} = \sum_{y_i = h_t(\mathbf{x}_i)} \frac{D_t(i) \exp(-\alpha_t)}{Z_t} = \frac{1}{2}$$

即

$$\sum_{i=1}^n D_{t+1}(i) \mathbb{1}_{[y_i \neq h_t(\mathbf{x}_i)]} = \sum_{i=1}^n D_{t+1}(i) \mathbb{1}_{[y_i = h_t(\mathbf{x}_i)]} = \frac{1}{2}$$

### 3.2 替换目标函数 (10pt)

1. 考虑如下的函数：(1)  $\phi_1(x) = \mathbb{1}_{x \geq 0}$ ；(2)  $\phi_2(x) = (1+x)^2$ ；(3)  $\phi_3(x) = \max\{0, 1+x\}$ ；(4)  $\phi_4(x) = \log_2(1+e^x)$ 。哪一个函数满足题面中对于  $\phi$  的假设条件？请求出使用该函数时， $D_t(i)$  的表达式。

显然只有  $\phi_4(x) = \log_2(1+e^x)$  是单调递增且处处可微的凸函数，且满足  $\forall x \geq 0, \phi(x) \geq 1$  且  $\forall x < 0, \phi(x) > 0$  的条件。

对于  $\phi_4(x) = \log_2(1+e^x)$ ，有

$$\phi'_4(x) = \frac{1}{\ln 2} \frac{e^x}{1+e^x} = \frac{1}{\ln 2} \frac{1}{1+e^{-x}}$$

故

$$D_t(i) = \frac{\phi'_4(-y_i f_t(\mathbf{x}_i))}{Z_t} = \frac{1}{Z_t \ln 2} \frac{1}{1+e^{y_i f_t(\mathbf{x}_i)}}$$

2. 上文的分析只确定了坐标下降中每次选择的最优坐标方向；对于前一问中选择的函数  $\phi$ ，请通过求解  $\frac{dL(\alpha + \beta e_t)}{d\beta} = 0$ ，进一步确定最优步长  $\beta$ 。（写出  $\beta$  应满足的方程并适当化简即可，无需求解方程）

目标函数为

$$L(\alpha) = \sum_{i=1}^n \log_2 \left( 1 + e^{-y_i f(\mathbf{x}_i)} \right)$$

其中

$$f(\mathbf{x}_i) = \sum_{t=1}^T \alpha_t h_t(\mathbf{x}_i)$$

引入步长  $\beta$  后，目标函数变为

$$L(\alpha + \beta e_t) = \sum_{i=1}^n \log_2 \left( 1 + e^{-y_i(f(\mathbf{x}_i) + \beta h_t(\mathbf{x}_i))} \right)$$

对  $\beta$  求导，有

$$\frac{dL(\alpha + \beta e_t)}{d\beta} = \frac{1}{\ln 2} \sum_{i=1}^n -\frac{y_i h_t(\mathbf{x}_i) e^{-y_i(f(\mathbf{x}_i) + \beta h_t(\mathbf{x}_i))}}{1 + e^{-y_i(f(\mathbf{x}_i) + \beta h_t(\mathbf{x}_i))}}$$

最优步长  $\beta$  应满足上式等于 0，适当化简后即为

$$\sum_{i=1}^n \frac{y_i h_t(\mathbf{x}_i)}{1 + e^{y_i(f(\mathbf{x}_i) + \beta h_t(\mathbf{x}_i))}} = 0$$

### 3.3 带未知标签的 Boosting 算法（附加题，10pt）

1. 用  $\epsilon_t^s$  和  $\alpha_t$  表示  $Z_t$ 。

$$\begin{aligned} Z_t &= \sum_{i=1}^n D_t(i) \exp(-\alpha_t y_i h_t(\mathbf{x}_i)) \\ &= \sum_{s=-1}^1 \sum_{i=1}^n D_t(i) \exp(-\alpha_t s) \mathbb{1}_{y_i h_t(\mathbf{x}_i)=s} \\ &= \sum_{s=-1}^1 \epsilon_t^s \exp(-\alpha_t s) \\ &= \epsilon_t^- e^{\alpha_t} + \epsilon_t^0 + \epsilon_t^+ e^{-\alpha_t} \end{aligned}$$

2. 计算  $F'(\bar{\alpha}_{t-1}, e_k)$ ，并指出第  $t$  步时弱分类器的优化目标。（结果用含  $\epsilon_t^s$  的表达式表示）

由于

$$\begin{aligned} F(\bar{\alpha}_{t-1}) &= \frac{1}{n} \sum_{i=1}^n e^{-y_i \sum_{j=1}^N \bar{\alpha}_{t-1,j} h_j(\mathbf{x}_i)} \\ F(\bar{\alpha}_{t-1} + \eta e_k) &= \frac{1}{n} \sum_{i=1}^n e^{-y_i \sum_{j=1}^N \bar{\alpha}_{t-1,j} h_j(\mathbf{x}_i) - \eta y_i h_k(\mathbf{x}_i)} \end{aligned}$$

故而有

$$\begin{aligned} F'(\bar{\alpha}_{t-1}, e_k) &= \lim_{\eta \rightarrow 0} \frac{F(\bar{\alpha}_{t-1} + \eta e_k) - F(\bar{\alpha}_{t-1})}{\eta} \\ &= -\frac{1}{n} \sum_{i=1}^n y_i h_k(\mathbf{x}_i) e^{-y_i \sum_{j=1}^N \bar{\alpha}_{t-1,j} h_j(\mathbf{x}_i)} \\ &= -\frac{1}{n} \sum_{i=1}^n y_i h_k(\mathbf{x}_i) \bar{D}_t(i) \bar{Z}_t \\ &= -\left[ \sum_{i=1}^n \bar{D}_t(i) \mathbb{1}_{y_i h_t(\mathbf{x}_i)=1} - \sum_{i=1}^n \bar{D}_t(i) \mathbb{1}_{y_i h_t(\mathbf{x}_i)=-1} \right] \frac{\bar{Z}_t}{n} \\ &= -[\epsilon_t^+ - \epsilon_t^-] \frac{\bar{Z}_t}{n} \end{aligned}$$

第  $t$  步时弱分类器的优化目标为在分布  $\bar{D}_t$  下，使得  $|\epsilon_t^+ - \epsilon_t^-|$  最小。

3. 解  $\frac{\partial F(\bar{\alpha}_{t-1} + \eta \mathbf{e}_k)}{\partial \eta} = 0$  , 并给出  $\alpha_t$  的更新公式。(结果用含  $\epsilon_t^s$  的表达式表示)

$$\begin{aligned}\frac{\partial F(\bar{\alpha}_{t-1} + \eta \mathbf{e}_k)}{\partial \eta} &= 0 \\ -\frac{1}{n} \sum_{i=1}^n y_i h_k(\mathbf{x}_i) e^{-y_i \sum_{j=1}^N \bar{\alpha}_{t-1,j} h_j(\mathbf{x}_i) - \eta y_i h_k(\mathbf{x}_i)} &= 0 \\ \sum_{i=1}^n y_i h_k(\mathbf{x}_i) \bar{D}_t(i) \bar{Z}_t e^{-\eta y_i h_k(\mathbf{x}_i)} &= 0 \\ \sum_{i=1}^n \bar{D}_t(i) y_i h_k(\mathbf{x}_i) e^{-\eta y_i h_k(\mathbf{x}_i)} &= 0 \\ \sum_{i=1}^n \bar{D}_t(i) \mathbb{1}_{y_i h_t(\mathbf{x}_i)=1} e^{-\eta} - \sum_{i=1}^n \bar{D}_t(i) \mathbb{1}_{y_i h_t(\mathbf{x}_i)=-1} e^{\eta} &= 0 \\ \epsilon_t^+ e^{-\eta} - \epsilon_t^- e^{\eta} &= 0\end{aligned}$$

故可以解得

$$\eta = \frac{1}{2} \log \frac{\epsilon_t^+}{\epsilon_t^-}$$

故而有

$$\alpha_t = \alpha_{t-1} + \frac{\mathbf{e}_k}{2} \log \frac{\epsilon_t^+}{\epsilon_t^-}$$

4. 在 AdaBoost 中, 我们证明了训练误差  $\hat{\mathcal{E}} = \frac{1}{n} \sum_{i=1}^n \mathbb{1}_{y_i f(\mathbf{x}_i) < 0} \leq \prod_{t=1}^T Z_t \leq \prod_{t=1}^T 2\sqrt{\epsilon_t(1-\epsilon_t)}$ 。

类似地, 请求出带未知标签的 Boosting 算法中, 用含  $\epsilon_t^s$  的表达式表示的训练误差上界, 并证明若每个弱学习器的误差满足  $\frac{\epsilon_t^+ - \epsilon_t^-}{\sqrt{1 - \epsilon_t^0}} \geq \gamma > 0$ , 则有:  $\frac{1}{n} \sum_{i=1}^n \mathbb{1}_{y_i f(\mathbf{x}_i) < 0} \leq \exp\left(-\frac{\gamma^2 T}{2}\right)$ , 其中,  $T$

是弱学习器的个数。

$$\begin{aligned}\hat{\mathcal{E}} &= \frac{1}{n} \sum_{i=1}^n \mathbb{1}_{y_i f(\mathbf{x}_i) < 0} \\ &\leq \frac{1}{n} \sum_{i=1}^n e^{-y_i f(\mathbf{x}_i)} \\ &= \frac{1}{n} \sum_{i=1}^n \left( \prod_{t=1}^T Z_t \right) D_{T+1}(i) \\ &= \prod_{t=1}^T Z_t\end{aligned}$$

上面已经证明  $Z_t = \epsilon_t^- e^{\alpha_t} + \epsilon_t^0 + \epsilon_t^+ e^{-\alpha_t}$ , 而通过  $\frac{\partial Z_t}{\partial \alpha_t} = \epsilon_t^- e^{\alpha_t} - \epsilon_t^+ e^{-\alpha_t} = 0$  可以解得

$\alpha_t = \frac{1}{2} \log \frac{\epsilon_t^+}{\epsilon_t^-}$ , 故而有

$$Z_t \leq \epsilon_t^0 + 2\sqrt{\epsilon_t^+ \epsilon_t^-}$$

故而有

$$\hat{\mathcal{E}} = \frac{1}{n} \sum_{i=1}^n \mathbb{1}_{y_i f(\mathbf{x}_i) < 0} \leq \prod_{t=1}^T Z_t \leq \prod_{t=1}^T \left( \epsilon_t^0 + 2\sqrt{\epsilon_t^+ \epsilon_t^-} \right)$$

注意到  $\epsilon_t^+ + \epsilon_t^- = 1 - \epsilon_t^0$ ，故  $Z_t \leq \epsilon_t^0 + 2\sqrt{\epsilon_t^+ \epsilon_t^-} = \epsilon_t^0 + \sqrt{(1 - \epsilon_t^0)^2 - (\epsilon_t^+ - \epsilon_t^-)^2}$ 。

故当  $\frac{\epsilon_t^+ - \epsilon_t^-}{\sqrt{1 - \epsilon_t^0}} \geq \gamma > 0$  即  $(\epsilon_t^+ - \epsilon_t^-)^2 \geq \gamma^2 (1 - \epsilon_t^0)$  时，有

$$\begin{aligned} \hat{\mathcal{E}} &= \frac{1}{n} \sum_{i=1}^n \mathbb{1}_{y_i f(\mathbf{x}_i) < 0} \\ &\leq \prod_{t=1}^T \left( \epsilon_t^0 + \sqrt{(1 - \epsilon_t^0)^2 - (\epsilon_t^+ - \epsilon_t^-)^2} \right) \\ &\leq \prod_{t=1}^T \left( \epsilon_t^0 + \sqrt{(1 - \epsilon_t^0)^2 - \gamma^2 (1 - \epsilon_t^0)} \right) \\ &\leq \prod_{t=1}^T \left( \epsilon_t^0 + (1 - \epsilon_t^0) \sqrt{1 - \gamma^2} \right) \\ &\leq \prod_{t=1}^T \sqrt{\epsilon_t^0 + (1 - \epsilon_t^0) (1 - \gamma^2)} \\ &= \prod_{t=1}^T \sqrt{1 - (1 - \epsilon_t^0) \gamma^2} \\ &\leq \prod_{t=1}^T \sqrt{\exp(-(1 - \epsilon_t^0) \gamma^2)} \\ &\leq \prod_{t=1}^T \sqrt{\exp(-\gamma^2)} \\ &= \exp\left(-\frac{\gamma^2 T}{2}\right) \end{aligned}$$

### 3.4 Gradient Boosting Machines (40pt)

1. 完成上述算法中的填空。

1. 令  $f_0(\mathbf{x}) = 0$ 。

2. For  $t = 1$  to  $T$ :

(a) 计算在各个数据点上的梯度  $\mathbf{g}_t = \left( \frac{\partial}{\partial \hat{\mathbf{y}}_i} \ell(\mathbf{y}_i, \hat{\mathbf{y}}_i) \Big|_{\hat{\mathbf{y}}_i = f_{t-1}(\mathbf{x}_i)} \right)_{i=1}^n$ 。

(b) 根据  $-\mathbf{g}_t$  拟合一个回归模型， $h_t = \arg \min_{h \in \mathcal{F}} \sum_{i=1}^n (h(\mathbf{x}_i) - \mathbf{g}_t)^2$ 。

(c) 选择合适的步长  $\alpha_t$ ，最简单的选择是固定步长  $\eta \in (0, 1]$ 。

(d) 更新模型， $f(\mathbf{x}) = f_{t-1}(\mathbf{x}) + \eta h_t(\mathbf{x})$ 。

2. 考虑回归问题, 假设损失函数  $\ell(\mathbf{y}, \hat{\mathbf{y}}) = \frac{1}{2}(\mathbf{y} - \hat{\mathbf{y}})^2$ 。直接给出第  $t$  轮迭代时的  $\mathbf{g}_t$  以及  $h_t$  的表达式。(使用  $f_{t-1}$  表达)。

$$\begin{aligned}\mathbf{g}_t &= \left. \frac{\partial}{\partial \hat{\mathbf{y}}} \frac{1}{2}(\mathbf{y} - \hat{\mathbf{y}})^2 \right|_{\hat{\mathbf{y}}=f_{t-1}(\mathbf{x})} \\ &= f_{t-1}(\mathbf{x}) - \mathbf{y} \\ h_t &= \arg \min_{h \in \mathcal{F}} \sum_{i=1}^n (h(\mathbf{x}_i) - \mathbf{g}_t)^2 \\ &= \arg \min_{h \in \mathcal{F}} \sum_{i=1}^n (h(\mathbf{x}_i) - f_{t-1}(\mathbf{x}_i) + \mathbf{y})^2\end{aligned}$$

3. 考虑二分类问题, 假设损失函数  $\ell(\mathbf{y}, \hat{\mathbf{y}}) = \ln(1 + e^{-\mathbf{y}\hat{\mathbf{y}}})$ 。直接给出第  $t$  轮迭代时的  $\mathbf{g}_t$  以及  $h_t$  的表达式。(使用  $f_{t-1}$  表达)。

$$\begin{aligned}\mathbf{g}_t &= \left. \frac{\partial}{\partial \hat{\mathbf{y}}} \ln(1 + e^{-\mathbf{y}\hat{\mathbf{y}}}) \right|_{\hat{\mathbf{y}}=f_{t-1}(\mathbf{x})} \\ &= \frac{-\mathbf{y}e^{-\mathbf{y}f_{t-1}(\mathbf{x})}}{1 + e^{-\mathbf{y}f_{t-1}(\mathbf{x})}} \\ &= -\mathbf{y}\sigma(-\mathbf{y}f_{t-1}(\mathbf{x})) \\ h_t &= \arg \min_{h \in \mathcal{F}} \sum_{i=1}^n (h(\mathbf{x}_i) - \mathbf{g}_t)^2 \\ &= \arg \min_{h \in \mathcal{F}} \sum_{i=1}^n (h(\mathbf{x}_i) + \mathbf{y}\sigma(-\mathbf{y}f_{t-1}(\mathbf{x}_i)))^2\end{aligned}$$

4. 完成 `boosting.py` 中 `GradientBoosting` 类的 `fit` 函数。

```
1 def fit(self, train_data, train_target):
2     """
3     Fit gradient boosting model
4     :param train_data: x
5     :param train_target: y
6     :return:
7     """
8     ft = np.zeros(train_data.shape[0]) # f_t(x)
9     train_target = train_target.squeeze()
10    self.ht = [] # sequence of h_t(x)
11    for t in range(self.T):
12        rgs = DecisionTreeRegressor(min_samples_split=self.min_sample,
13        max_depth=self.max_depth)
14        # Your code goes here (~4 lines)
15        # TODO 3.5.4
16        gradient = self.gradient_func(train_target, ft)
17        rgs.fit(train_data, -gradient)
18        self.ht.append(rgs)
19        ft += self.learning_rate * rgs.predict(train_data)
20    return self
```

5. 完成 `boosting.py` 中 `GradientBoosting` 类的 `predict` 函数。

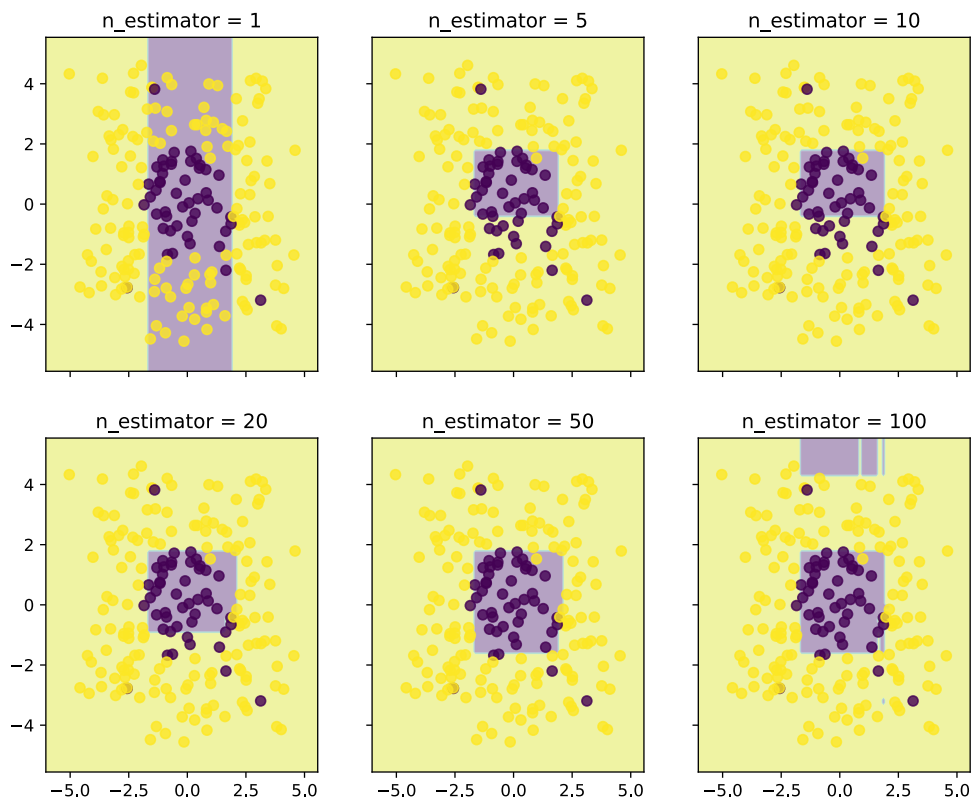
```
1 def predict(self, test_data):
2     # Your code goes here (~6 lines)
3     # TODO 3.5.5
4     ft = np.zeros(test_data.shape[0])
5     for rgs in self.ht:
6         ft += self.learning_rate * rgs.predict(test_data)
7     return ft
```

6. 完成 `boosting.py` 中函数 `gradient_logistic`。

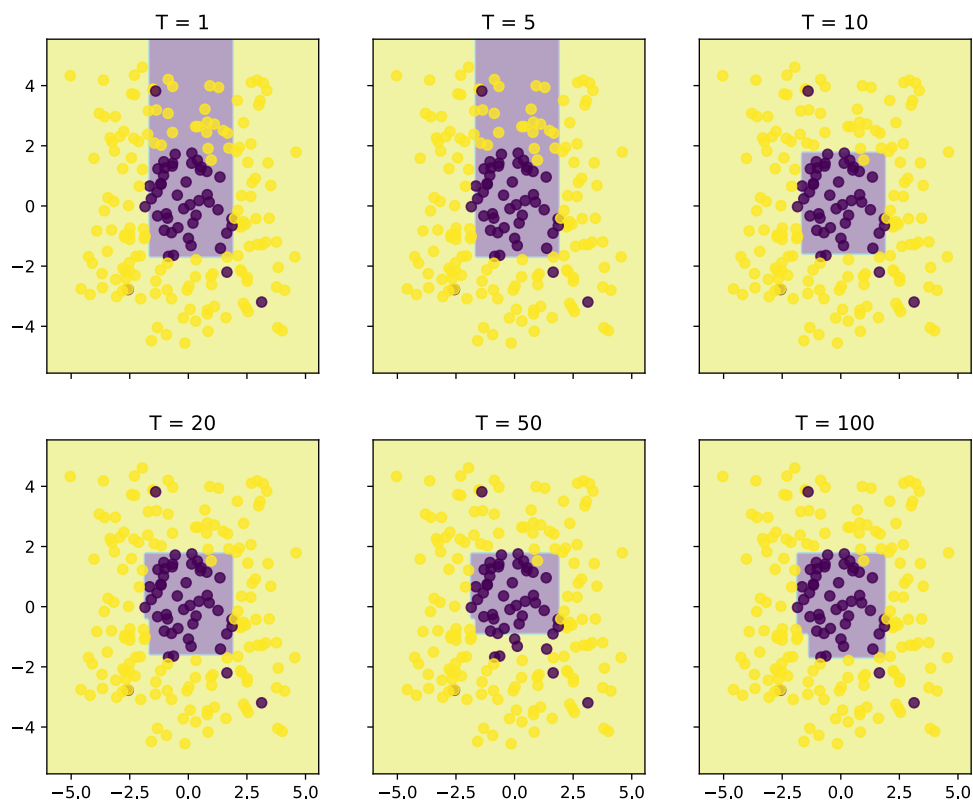
```
1 def gradient_logistic(train_target, train_predict):
2     """
3     compute g_t in 3.5.3
4     """
5     # Your code goes here (~3 lines)
6     # TODO 3.5.6
7     sigmoid = 1 / (1 + np.exp(train_target * train_predict))
8     return -train_target * sigmoid
```

7. 运行 `boosting.py`，在实验文档中记录 GBM 在不同数据集上运行的结果，包括 (a) `GBM_l2.pdf`，使用 L2 loss 在二分类问题上的结果。(b) `GBM_logistic.pdf`，使用 logistic loss 在二分类问题上的结果。(c) `GBM_regression.pdf`，使用 L2 loss 在回归问题上的结果。并简要描述实验现象（例如超参数对于 GBM 的影响、损失函数对于 GBM 的影响等）。

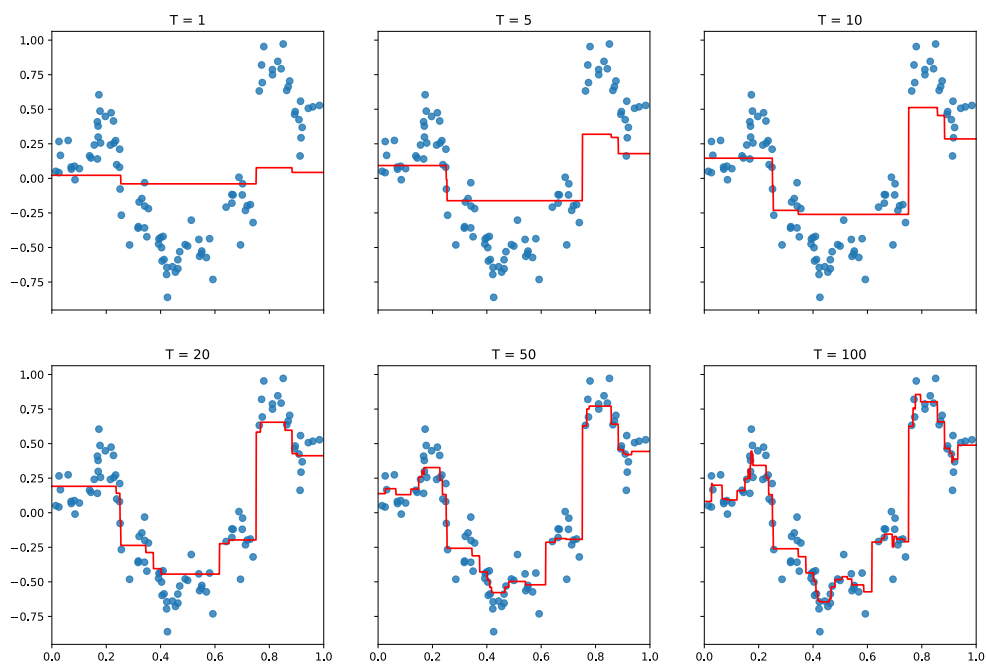
`GBM_l2.pdf` 的结果如下



`GBM_logistic.pdf` 的结果如下



GBM\_regression.pdf 的结果如下



可以观察到，对于 L2 loss 在二分类问题上的结果，当 `n_estimator` 较小时，存在欠拟合现象；随着 `n_estimator` 的增加，模型的拟合效果逐渐提升，但是当 `n_estimator` 较大时，存在过拟合现象。

对于 logistic loss 在二分类问题上的结果，当  $T$  较小时，存在欠拟合现象；随着  $T$  的增加，模型的拟合效果逐渐提升，且整体上优于 L2 loss。

对于 L2 loss 在回归问题上的结果，当  $T$  较小时，存在欠拟合现象；随着  $T$  的增加，模型的拟合效果逐渐提升，但是当  $T$  较大时，存在少量的过拟合现象，但整体效果上优于决策树的结果。

---

成绩：100

评语：2.1.2 -3; 3.3（附加题） 10