

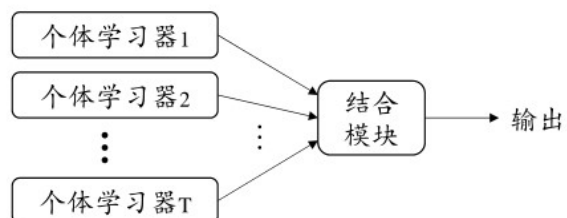
第八章：集成学习

集成学习

- 个体与集成
- Boosting
 - Adaboost
- Bagging与随机森林
- 结合策略
 - 平均法
 - 投票法
 - 学习法
- 多样性
 - 误差-分歧分解
 - 多样性度量
 - 多样性扰动

个体与集成

- 集成学习(ensemble learning)通过构建并结合多个学习器来提升性能



个体与集成

- 考虑一个简单的例子，在二分类问题中，假定3个分类器在三个样本中的表现如下图所示，其中√表示分类正确，X号表示分类错误，集成的结果通过投票产生。

| 测试例1 | 测试例2 | 测试例3 | 测试例1 | 测试例2 | 测试例3 | 测试例1 | 测试例2 | 测试例3 |
|------------|------|------|------------|-------|------|------------|------|------|
| h_1 | √ | √ | × | h_1 | √ | √ | × | × |
| h_2 | × | √ | √ | h_2 | √ | √ | × | × |
| h_3 | √ | × | √ | h_3 | √ | √ | × | × |
| 集群 | √ | √ | √ | 集群 | √ | √ | × | × |
| (a) 集群提升性能 | | | (b) 集群不起作用 | | | (c) 集群起负作用 | | |

- 集成个体应：好而不同

个体与集成 - 简单分析

- 考虑二分类问题，假设基分类器的错误率为：

$$P(h_i(\mathbf{x}) \neq f(\mathbf{x})) = \epsilon$$

- 假设集成通过简单投票法结合 T 个分类器，若有超过半数的基分类器正确则分类就正确

$$H(\mathbf{x}) = \text{sign}\left(\sum_{i=1}^T h_i(\mathbf{x})\right)$$

个体与集成 - 简单分析

- 假设基分类器的错误率相互独立，则由Hoeffding不等式可得集成的错误率为：

$$\begin{aligned} P(H(\mathbf{x}) \neq f(\mathbf{x})) &= \sum_{k=0}^{\lfloor T/2 \rfloor} \binom{T}{k} (1-\epsilon)^k \epsilon^{T-k} \\ &\leq \exp\left(-\frac{1}{2}T(1-2\epsilon)^2\right) \end{aligned}$$

- 上式显示，在一定条件下，随着集成分类器数目的增加，集成的错误率将指数级下降，最终趋向于0

个体与集成 - 简单分析

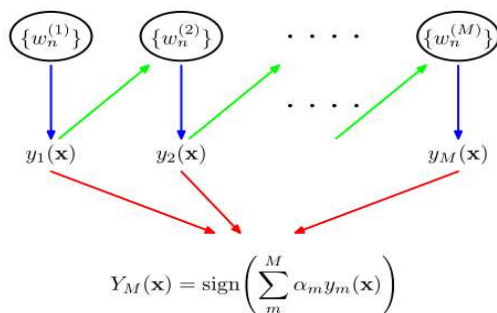
- 上面的分析有一个关键假设：基学习器的误差相互独立
- 现实任务中，个体学习器是为解决同一个问题训练出来的，显然不可能互相独立
- 事实上，个体学习器的“准确性”和“多样性”本身就存在冲突
- 如何产生“好而不同”的个体学习器是集成学习研究的核心
- 集成学习大致可分为两大类

集成学习

- 个体与集成
- Boosting
 - Adaboost
- Bagging与随机森林
- 结合策略
 - 平均法
 - 投票法
 - 学习法
- 多样性
 - 误差-分歧分解
 - 多样性度量
 - 多样性扰动

Boosting

- 个体学习器存在强依赖关系,
- 串行生成
- 每次调整训练数据的样本分布



Boosting - Boosting算法

Input: Sample distribution \mathcal{D} ;
Base learning algorithm \mathcal{L} ;
Number of learning rounds T .

Process:

1. $\mathcal{D}_1 = \mathcal{D}$. % Initialize distribution
2. **for** $t = 1, \dots, T$:
3. $h_t = \mathcal{L}(\mathcal{D}_t)$; % Train a weak learner from distribution \mathcal{D}_t
4. $\epsilon_t = P_{\mathbf{x} \sim \mathcal{D}_t}(h_t(\mathbf{x}) \neq f(\mathbf{x}))$; % Evaluate the error of h_t
5. $\mathcal{D}_{t+1} = \text{Adjust_Distribution}(\mathcal{D}_t, \epsilon_t)$
6. **end**

Output: $H(\mathbf{x}) = \text{Combine_Outputs}(\{h_1(\mathbf{x}), \dots, h_t(\mathbf{x})\})$

□ Boosting族算法最著名的代表是AdaBoost

Boosting – AdaBoost算法

输入: 训练集 $D = \{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_m, y_m)\}$;
基学习算法 \mathcal{L} ;
训练轮数 T .

过程:

- 1: $\mathcal{D}_1(\mathbf{x}) = 1/m$.
- 2: **for** $t = 1, 2, \dots, T$ **do**
- 3: $h_t = \mathcal{L}(D, \mathcal{D}_t)$;
- 4: $\epsilon_t = P_{\mathbf{x} \sim \mathcal{D}_t}(h_t(\mathbf{x}) \neq f(\mathbf{x}))$;
- 5: **if** $\epsilon_t > 0.5$ **then break**
- 6: $\alpha_t = \frac{1}{2} \ln \left(\frac{1-\epsilon_t}{\epsilon_t} \right)$;
- 7: $\mathcal{D}_{t+1}(\mathbf{x}) = \frac{\mathcal{D}_t(\mathbf{x})}{Z_t} \times \begin{cases} \exp(-\alpha_t), & \text{if } h_t(\mathbf{x}) = f(\mathbf{x}) \\ \exp(\alpha_t), & \text{if } h_t(\mathbf{x}) \neq f(\mathbf{x}) \end{cases}$

$$= \frac{\mathcal{D}_t(\mathbf{x}) \exp(-\alpha_t f(\mathbf{x}) h_t(\mathbf{x}))}{Z_t}$$

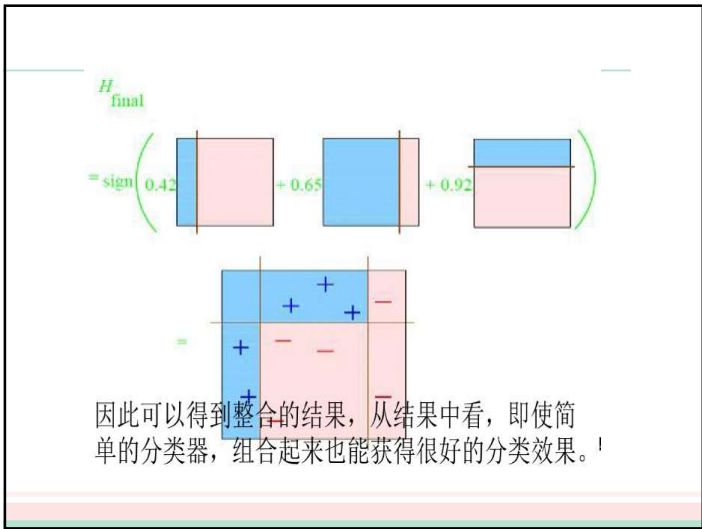
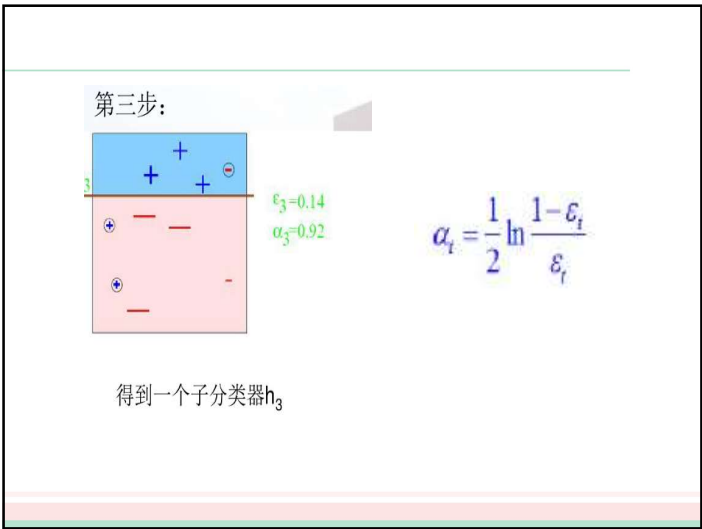
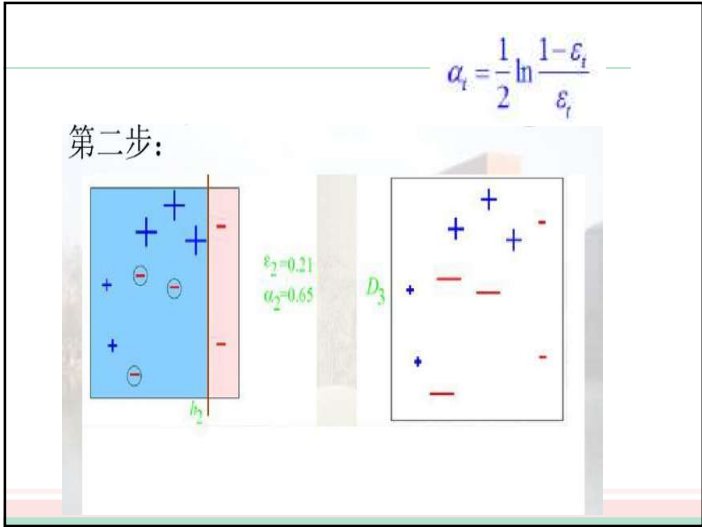
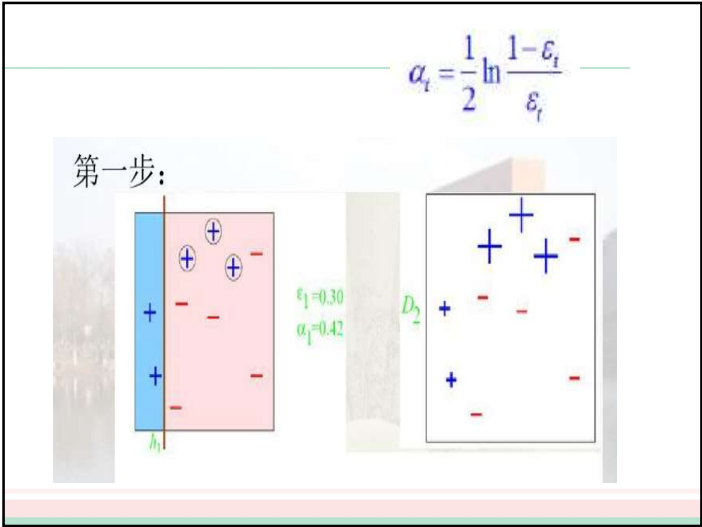
8: **end for**

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

Boosting

□ **AdaBoost** (Adaptive boosting) 算法:

- 初始化训练数据的权值分布。如果有N个样本, 则每一个训练样本最开始时都被赋予相同的权值: $1/N$ 。
- 训练弱分类器。具体训练过程中, 如果某个样本点已经被准确地分类, 那么在构造下一个训练集中, 它的权值就被降低; 相反, 如果某个样本点没有被准确地分类, 那么它的权值就得到提高。然后, 权值更新过的样本集被用于训练下一个分类器, 整个训练过程如此迭代地进行下去。
- 将各个训练得到的弱分类器组合成强分类器。各个弱分类器的训练过程结束后, 加大分类误差率小的弱分类器的权重, 使其在最终的分函数中起着较大的决定作用, 而降低分类误差率大的弱分类器的权重, 使其在最终的分函数中起着较小的决定作用。换言之, 误差率低的弱分类器在最终分类器中占的权重较大, 否则较小。



Boosting – AdaBoost推导

- 基学习器的线性组合

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

- 最小化指数损失函数

$$\ell_{\text{exp}}(H \mid \mathcal{D}) = \mathbb{E}_{\mathbf{x} \sim \mathcal{D}}[e^{-f(\mathbf{x})H(\mathbf{x})}]$$

Boosting – AdaBoost推导

- 若 $H(\mathbf{x})$ 能令指数损失函数最小化，则上式对 $H(\mathbf{x})$ 的偏导值为0，即

$$\frac{\partial \ell_{\text{exp}}(H \mid \mathcal{D})}{\partial H(\mathbf{x})} = -e^{-H(\mathbf{x})} P(f(\mathbf{x}) = 1 \mid \mathbf{x}) + e^{H(\mathbf{x})} P(f(\mathbf{x}) = -1 \mid \mathbf{x})$$

$$\begin{aligned} \text{sign}(H(\mathbf{x})) &= \text{sign}\left(\frac{1}{2} \ln \frac{P(f(\mathbf{x}) = 1 \mid \mathbf{x})}{P(f(\mathbf{x}) = -1 \mid \mathbf{x})}\right) & H(\mathbf{x}) &= \frac{1}{2} \ln \frac{P(f(\mathbf{x}) = 1 \mid \mathbf{x})}{P(f(\mathbf{x}) = -1 \mid \mathbf{x})} \\ &= \begin{cases} 1, & P(f(\mathbf{x}) = 1 \mid \mathbf{x}) > P(f(\mathbf{x}) = -1 \mid \mathbf{x}) \\ -1, & P(f(\mathbf{x}) = 1 \mid \mathbf{x}) < P(f(\mathbf{x}) = -1 \mid \mathbf{x}) \end{cases} \\ &= \arg \max_{y \in \{-1, 1\}} P(f(\mathbf{x}) = y \mid \mathbf{x}), \end{aligned}$$

$\text{sign}(H(\mathbf{x}))$ 达到了贝叶斯最优错误率，

说明指数损失函数是分类任务原来0/1

损失函数的一致的替代函数。

Boosting – AdaBoost推导

- 当基分类器 h_t 基于分布 D_t 产生后，该基分类器的权重 α_t 应使得 $\alpha_t h_t$ 最小化指数损失函数

$$\begin{aligned} \ell_{\text{exp}}(\alpha_t h_t \mid \mathcal{D}_t) &= \mathbb{E}_{\mathbf{x} \sim \mathcal{D}_t} [e^{-f(\mathbf{x})\alpha_t h_t(\mathbf{x})}] \\ &= \mathbb{E}_{\mathbf{x} \sim \mathcal{D}_t} [e^{-\alpha_t} \mathbb{I}(f(\mathbf{x}) = h_t(\mathbf{x})) + e^{\alpha_t} \mathbb{I}(f(\mathbf{x}) \neq h_t(\mathbf{x}))] \\ &= e^{-\alpha_t} P_{\mathbf{x} \sim \mathcal{D}_t}(f(\mathbf{x}) = h_t(\mathbf{x})) + e^{\alpha_t} P_{\mathbf{x} \sim \mathcal{D}_t}(f(\mathbf{x}) \neq h_t(\mathbf{x})) \\ &= e^{-\alpha_t} (1 - \epsilon_t) + e^{\alpha_t} \epsilon_t, \quad \epsilon_t = P_{\mathbf{x} \sim \mathcal{D}_t}(h_t(\mathbf{x}) \neq f(\mathbf{x})) \end{aligned}$$

- 令指数损失函数的导数为0，即

$$\frac{\partial \ell_{\text{exp}}(\alpha_t h_t \mid \mathcal{D}_t)}{\partial \alpha_t} = -e^{-\alpha_t} (1 - \epsilon_t) + e^{\alpha_t} \epsilon_t \quad \alpha_t = \frac{1}{2} \ln \left(\frac{1 - \epsilon_t}{\epsilon_t} \right)$$

Boosting – AdaBoost推导

- 在获得 H_{t-1} 之后的样本分布进行调整，使得下一轮的基学习器 h_t 能纠正 H_{t-1} 的一些错误。理想的 h_t 能纠正全部错误

$$\begin{aligned} \ell_{\text{exp}}(H_{t-1} + h_t \mid \mathcal{D}) &= \mathbb{E}_{\mathbf{x} \sim \mathcal{D}} [e^{-f(\mathbf{x})(H_{t-1}(\mathbf{x}) + h_t(\mathbf{x}))}] \\ &= \mathbb{E}_{\mathbf{x} \sim \mathcal{D}} [e^{-f(\mathbf{x})H_{t-1}(\mathbf{x})} e^{-f(\mathbf{x})h_t(\mathbf{x})}] \end{aligned}$$

- 泰勒展开近似为

$$\begin{aligned} \ell_{\text{exp}}(H_{t-1} + h_t \mid \mathcal{D}) &\simeq \mathbb{E}_{\mathbf{x} \sim \mathcal{D}} \left[e^{-f(\mathbf{x})H_{t-1}(\mathbf{x})} \left(1 - f(\mathbf{x})h_t(\mathbf{x}) + \frac{f^2(\mathbf{x})h_t^2(\mathbf{x})}{2} \right) \right] \\ &= \mathbb{E}_{\mathbf{x} \sim \mathcal{D}} \left[e^{-f(\mathbf{x})H_{t-1}(\mathbf{x})} \left(1 - f(\mathbf{x})h_t(\mathbf{x}) + \frac{1}{2} \right) \right] \end{aligned}$$

Boosting – AdaBoost推导

□ 于是，理想的基学习器：

$$\begin{aligned} h_t(\mathbf{x}) &= \arg \min_h \ell_{\exp}(H_{t-1} + h \mid \mathcal{D}) \\ &= \arg \min_h \mathbb{E}_{\mathbf{x} \sim \mathcal{D}} \left[e^{-f(\mathbf{x})H_{t-1}(\mathbf{x})} \left(1 - f(\mathbf{x})h(\mathbf{x}) + \frac{1}{2} \right) \right] \\ &= \arg \max_h \mathbb{E}_{\mathbf{x} \sim \mathcal{D}} \left[e^{-f(\mathbf{x})H_{t-1}(\mathbf{x})} f(\mathbf{x})h(\mathbf{x}) \right] \\ &= \arg \max_h \mathbb{E}_{\mathbf{x} \sim \mathcal{D}} \left[\frac{e^{-f(\mathbf{x})H_{t-1}(\mathbf{x})}}{\mathbb{E}_{\mathbf{x} \sim \mathcal{D}}[e^{-f(\mathbf{x})H_{t-1}(\mathbf{x})}]} f(\mathbf{x})h(\mathbf{x}) \right], \end{aligned}$$

□ 注意到 $\mathbb{E}_{\mathbf{x} \sim \mathcal{D}}[e^{-f(\mathbf{x})H_{t-1}(\mathbf{x})}]$ 是一个常数，令 \mathcal{D}_t 表示一个分布：

$$\mathcal{D}_t(\mathbf{x}) = \frac{\mathcal{D}(\mathbf{x})e^{-f(\mathbf{x})H_{t-1}(\mathbf{x})}}{\mathbb{E}_{\mathbf{x} \sim \mathcal{D}}[e^{-f(\mathbf{x})H_{t-1}(\mathbf{x})}]}$$

Boosting – AdaBoost推导

□ 根据数学期望的定义，这等价于令：

$$\begin{aligned} h_t(\mathbf{x}) &= \arg \max_h \mathbb{E}_{\mathbf{x} \sim \mathcal{D}} \left[\frac{e^{-f(\mathbf{x})H_{t-1}(\mathbf{x})}}{\mathbb{E}_{\mathbf{x} \sim \mathcal{D}}[e^{-f(\mathbf{x})H_{t-1}(\mathbf{x})}]} f(\mathbf{x})h(\mathbf{x}) \right] \\ &= \arg \max_h \mathbb{E}_{\mathbf{x} \sim \mathcal{D}_t} [f(\mathbf{x})h(\mathbf{x})]. \end{aligned}$$

□ 由 $f(\mathbf{x}), h(\mathbf{x}) \in \{-1, +1\}$ 有：

$$f(\mathbf{x})h(\mathbf{x}) = 1 - 2\mathbb{I}(f(\mathbf{x}) \neq h(\mathbf{x}))$$

Boosting – AdaBoost推导

□ 则理想的基学习器

$$h_t(\mathbf{x}) = \arg \min_h \mathbb{E}_{\mathbf{x} \sim \mathcal{D}_t} [\mathbb{I}(f(\mathbf{x}) \neq h(\mathbf{x}))]$$

□ 最终的样本分布更新公式

$$\begin{aligned} \mathcal{D}_{t+1}(\mathbf{x}) &= \frac{\mathcal{D}(\mathbf{x})e^{-f(\mathbf{x})H_t(\mathbf{x})}}{\mathbb{E}_{\mathbf{x} \sim \mathcal{D}}[e^{-f(\mathbf{x})H_t(\mathbf{x})}]} \\ &= \frac{\mathcal{D}(\mathbf{x})e^{-f(\mathbf{x})H_{t-1}(\mathbf{x})}e^{-f(\mathbf{x})\alpha_t h_t(\mathbf{x})}}{\mathbb{E}_{\mathbf{x} \sim \mathcal{D}}[e^{-f(\mathbf{x})H_t(\mathbf{x})}]} \\ &= \mathcal{D}_t(\mathbf{x}) \cdot e^{-f(\mathbf{x})\alpha_t h_t(\mathbf{x})} \frac{\mathbb{E}_{\mathbf{x} \sim \mathcal{D}}[e^{-f(\mathbf{x})H_{t-1}(\mathbf{x})}]}{\mathbb{E}_{\mathbf{x} \sim \mathcal{D}}[e^{-f(\mathbf{x})H_t(\mathbf{x})}]} \end{aligned}$$

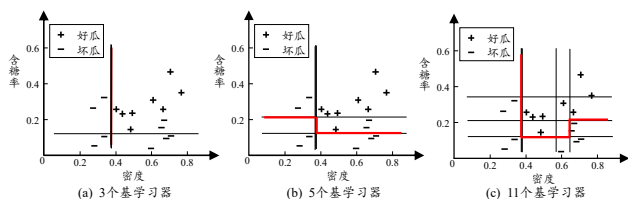
Boosting – AdaBoost注意事项

□ 数据分布的学习

- 重赋权法
- 重采样法

□ 重启，避免训练过程过早停止

Boosting – AdaBoost实验



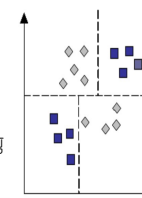
- 从偏差-方差的角度：降低偏差，可对泛化性能相当弱的学习器构造出很强的集成

Boosting+CART+任意损失函数=GBDT

CART假设决策树是二叉树，内部结点特征的取值为“是”和“否”，左分支是取值为“是”的分支，右分支是取值为“否”的分支。这样的决策树等价于递归地二分每个特征，将输入空间即特征空间划分为有限个单元，并在这些单元上确定预测的概率分布，也就是在输入给定的条件下输出的条件概率分布。

CART算法由以下两步组成：

- (1) 决策树的生成：基于训练数据集生成决策树，生成的决策树要尽量大（大是为了更好地泛化）；
- (2) 决策树剪枝：用验证数据集对已生成的树进行剪枝并选择最优子树，这时损失函数最小作为剪枝的标准。



GBDT模型表示为决策树的加法模型

$$f_M(x) = \sum_{m=1}^M T(x; \Theta_m)$$

其中， $T(x; \Theta_m)$ 表示决策树； m 为决策树的参数； M 为树的个数。

采用
$$f_m(x) = f_{m-1}(x) + T(x; \Theta_m)$$

要优化的损失函数

经验风险极小化确定下一棵树的参数

$$\hat{\Theta}_m = \arg \min_{\Theta_m} \sum_{i=1}^N L(y_i, f_{m-1}(x_i) + T(x_i; \Theta_m))$$

则上面损失函数变为：

$$\begin{aligned} L(y, f_{m-1}(x) + T(x; \Theta_m)) \\ &= [y - f_{m-1}(x) - T(x; \Theta_m)]^2 \\ &= [r - T(x; \Theta_m)]^2 \end{aligned}$$

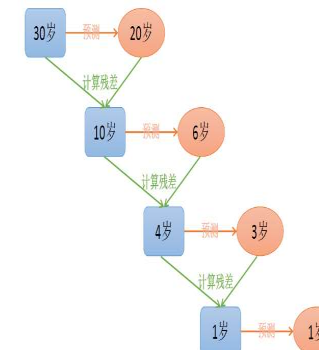
残差

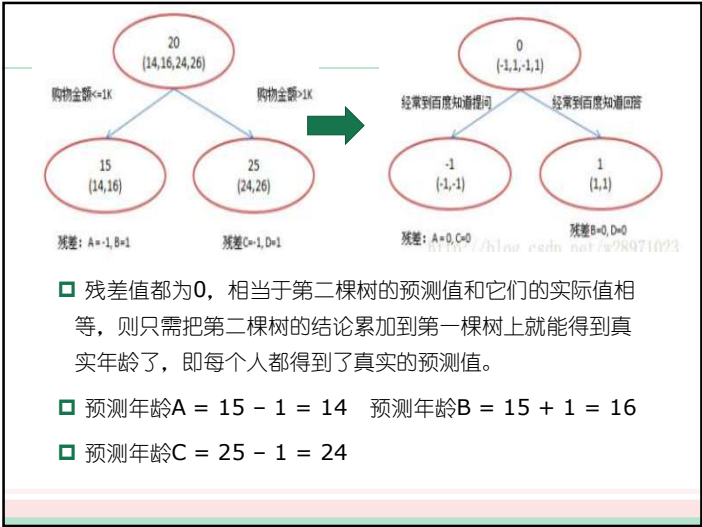
GBDT 在模型训练的时候，是要求模型预测的样本损失尽可能的小。

残差的意思就是：

- A 的实际值 - A 的预测值 = A 的残差

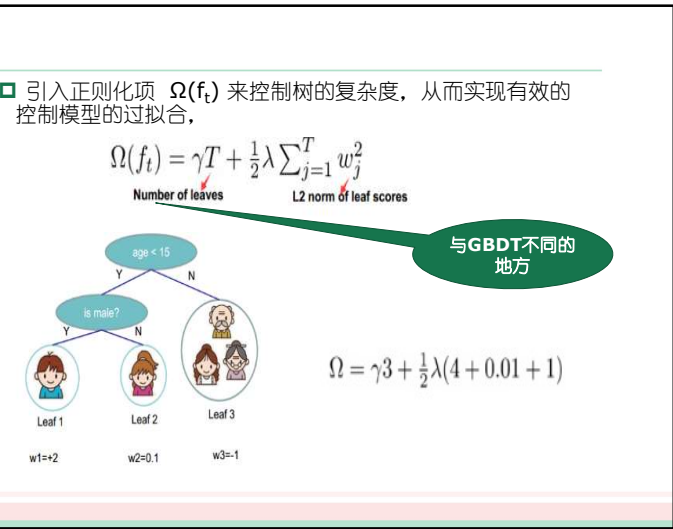
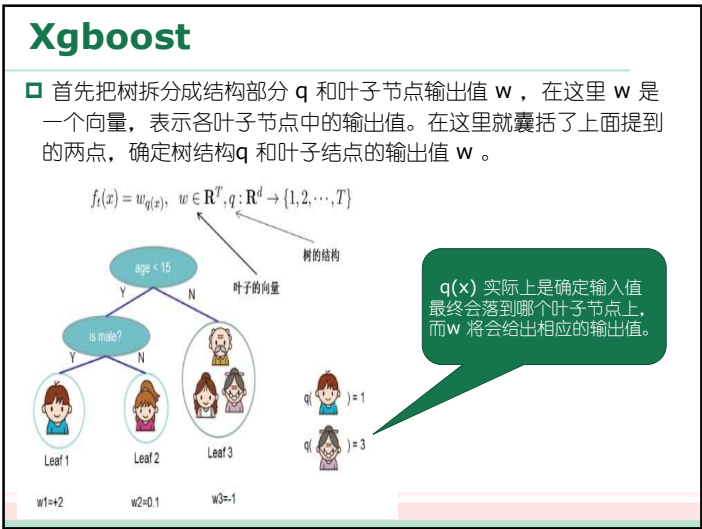
GBDT每一次的计算都是为了减少上一次的残差，进而在残差减少（负梯度）的方向上建立一个新的模型



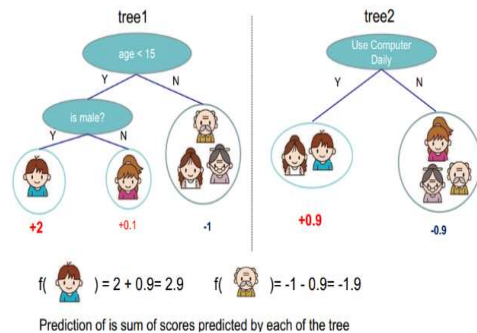


GBDT优缺点

- GBDT主要的优点有：
 - 可以灵活处理各种类型的数据，包括连续值和离散值。
 - 在相对少的调参时间情况下，预测的准确率也可以比较高。这个是相对SVM来说的。
 - 使用一些健壮的损失函数，对异常值的鲁棒性非常强。比如 Huber 损失函数和Quantile损失函数。
- GBDT的主要缺点有：
 - 由于弱学习器之间存在依赖关系，难以并行训练数据。不过可以通过自采样的SGBT来达到部分并行。



❑ XGBoost中的Boosting Tree模型



集成学习

- ❑ 个体与集成
- ❑ Boosting
 - Adaboost
- ❑ Bagging与随机森林
- ❑ 结合策略
 - 平均法
 - 投票法
 - 学习法
- ❑ 多样性
 - 误差-分歧分解
 - 多样性度量
 - 多样性扰动

Bagging与随机森林

- ❑ 个体学习器不存在强依赖关系
- ❑ 并行化生成
- ❑ 自助采样法

Bagging与随机森林 - Bagging算法

输入: 训练集 $D = \{(x_1, y_1), (x_2, y_2), \dots, (x_m, y_m)\}$;
基学习算法 \mathcal{L} ;
训练轮数 T .

过程:

- 1: for $t = 1, 2, \dots, T$ do
- 2: $h_t = \mathcal{L}(D, \mathcal{D}_{bs})$
- 3: end for

输出: $H(x) = \arg \max_{y \in \mathcal{Y}} \sum_{t=1}^T \mathbb{I}(h_t(x) = y)$

Bagging与随机森林 - Bagging算法特点

- 时间复杂度低
 - 假定基学习器的计算复杂度为 $O(m)$ ，采样与投票/平均过程的复杂度为 $O(s)$ ，则bagging的复杂度大致为 $T(O(m)+O(s))$
 - 由于 $O(s)$ 很小且 T 是一个不大的常数
 - 因此训练一个bagging集成与直接使用基学习器的复杂度同阶
- 可使用包外估计

Bagging与随机森林 - 包外估计

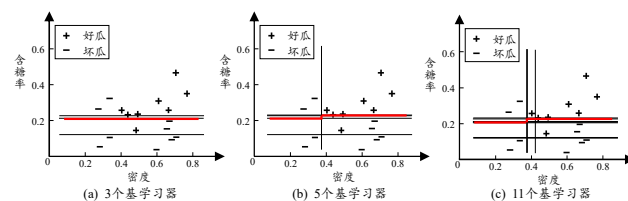
- $H^{oob}(x)$ 表示对样本 x 的包外预测，即仅考虑那些未使用样本 x 训练的基学习器在 x 上的预测

$$H^{oob}(x) = \arg \max_{y \in \mathcal{Y}} \sum_{t=1}^T \mathbb{I}(h_t(x) = y) \cdot \mathbb{I}(x \notin D_t)$$

- Bagging泛化误差的包外估计为：

$$\epsilon^{oob} = \frac{1}{|D|} \sum_{(x,y) \in D} \mathbb{I}(H^{oob}(x) \neq y)$$

Bagging与随机森林- Bagging实验

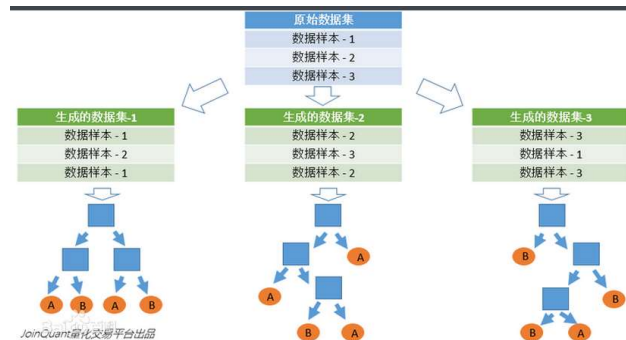


- 从偏差-方差的角度：降低方差，在不剪枝的决策树、神经网络等易受样本影响的学习器上效果更好

Bagging与随机森林-随机森林

- 随机森林(Random Forest, 简称RF)是bagging的一个扩展变种
- 采样的随机性
- 属性选择的随机性

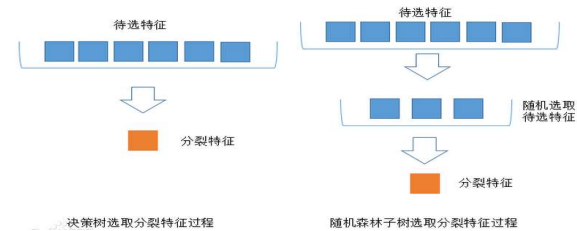
数据随机抽取



41

特征随机抽取

□ 与数据集的随机选取类似，随机森林中的子树的每一个分裂过程并未用到所有的待选特征，而是从所有的待选特征中随机选取一定的特征，之后在随机选取的特征中选取最优的特征。这样能够使得随机森林中的决策树都能够彼此不同，提升系统的多样性，从而提升分类性能。



42

□ 每棵树的按照如下规则生成：

- 如果训练集大小为 N ，对于每棵树而言，随机且有放回地从训练集中的抽取 N 个训练样本（这种采样方式称为bootstrap sample方法），作为该树的训练集；
- 如果每个样本的特征维度为 M ，指定一个常数 $m \ll M$ ，随机地从 M 个特征中选取 m 个特征子集，每次树进行分裂时，从这 m 个特征中选择最优的；
- 每棵树都尽最大程度的生长，并且没有剪枝过程。

43

□ 随机森林分类效果（错误率）与两个因素有关：

- 森林中任意两棵树的相关性：相关性越大，错误率越大；
- 森林中每棵树的分类能力：每棵树的分类能力越强，整个森林的错误率越低。

减小特征选择个数 m ，树的相关性和分类能力也会相应的降低；增大 m ，两者也会随之增大。所以关键问题是如何选择最优的 m （或者是范围），这也是随机森林唯一的一个参数。

44

```

❑ from sklearn.ensemble import RandomForestClassifier
❑ # 开始利用机器学习算法计算, 括号里面的n_estimators就是森林中包含的树的数目啦
❑ clf = RandomForestClassifier(n_estimators=10)
❑ # 训练的代码clf.fit(x_train, y_train)
❑ #
❑ 得到测试结果的代码prediction = clf.predict(x_test)
❑ # 看看预测对了没
❑ print(prediction == y_test)
❑ print('all done')
❑
❑ 输出:
❑ [ True]all done

```

45

Bagging与随机森林 - 随机森林算法

❑ 随机森林算法

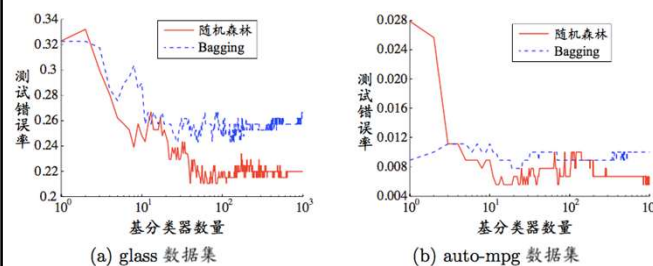
Input: Data set $D = \{(x_1, y_1), (x_2, y_2), \dots, (x_m, y_m)\}$;
Feature subset size K .

Process:

1. $N \leftarrow$ create a tree node based on D ;
2. **if all instances in the same class then return** N
3. $\mathcal{F} \leftarrow$ the set of features that can be split further;
4. **if \mathcal{F} is empty then return** N
5. $\tilde{\mathcal{F}} \leftarrow$ select K features from \mathcal{F} randomly;
6. $N.f \leftarrow$ the feature which has the best split point in $\tilde{\mathcal{F}}$;
7. $N.p \leftarrow$ the best split point on $N.f$;
8. $D_l \leftarrow$ subset of D with values on $N.f$ smaller than $N.p$;
9. $D_r \leftarrow$ subset of D with values on $N.f$ no smaller than $N.p$;
10. $N_l \leftarrow$ call the process with parameters (D_l, K) ;
11. $N_r \leftarrow$ call the process with parameters (D_r, K) ;
12. **return** N

Output: A random decision tree

Bagging与随机森林 - 随机森林实验

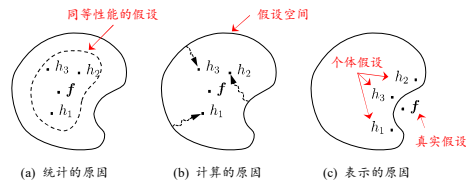


集成学习

- ❑ 个体与集成
- ❑ Boosting
 - Adaboost
- ❑ Bagging与随机森林
- ❑ 结合策略
 - 平均法
 - 投票法
 - 学习法
- ❑ 多样性
 - 误差-分歧分解
 - 多样性度量
 - 多样性扰动

结合策略

- 学习器的组合可以从三个方面带来好处



结合策略 - 平均法

- 简单平均法

$$H(\mathbf{x}) = \frac{1}{T} \sum_{i=1}^T h_i(\mathbf{x}).$$

- 加权平均法

$$H(\mathbf{x}) = \sum_{i=1}^T w_i h_i(\mathbf{x}), \quad w_i \geq 0 \text{ and } \sum_{i=1}^T w_i = 1.$$

结合策略 - 平均法

- 简单平均法是加权平均法的特例
- 加权平均法在二十世纪五十年代被广泛使用
- 集成学习中的各种结合方法都可以看成是加权平均法的变种或特例
- 加权平均法可认为是集成学习研究的基本出发点
- 加权平均法未必一定优于简单平均法

结合策略 - 投票法

- 绝对多数投票法 (majority voting)

$$H(\mathbf{x}) = \begin{cases} c_j & \text{if } \sum_{i=1}^T h_i^j(\mathbf{x}) > \frac{1}{2} \sum_{k=1}^l \sum_{i=1}^T h_i^k(\mathbf{x}) \\ \text{rejection} & \text{otherwise} \end{cases}$$

- 相对多数投票法 (plurality voting)

$$H(\mathbf{x}) = c_{\arg \max_j \sum_{i=1}^T h_i^j(\mathbf{x})}$$

- 加权投票法 (weighted voting)

$$H(\mathbf{x}) = c_{\arg \max_j \sum_{i=1}^T w_i h_i^j(\mathbf{x})}$$

结合策略 - 学习法

Stacking是学习法的典型代表

Input: Data set $D = \{(x_1, y_1), (x_2, y_2), \dots, (x_m, y_m)\}$;
First-level learning algorithms $\mathcal{L}_1, \dots, \mathcal{L}_T$;
Second-level learning algorithm \mathcal{L} .

Process:

```

1. for  $t = 1, \dots, T$ : % Train a first-level learner by applying the
2.    $h_t = \mathcal{L}_t(D)$ ; % first-level learning algorithm  $\mathcal{L}_t$ 
3. end
4.  $D' = \emptyset$ ; % Generate a new data set
5. for  $i = 1, \dots, m$ :
6.   for  $t = 1, \dots, T$ :
7.      $z_{it} = h_t(x_i)$ ;
8.   end
9.    $D' = D' \cup \{(z_{i1}, \dots, z_{iT}), y_i\}$ ;
10. end
11.  $h' = \mathcal{L}(D')$ ; % Train the second-level learner  $h'$  by applying
    % the second-level learning algorithm  $\mathcal{L}$  to the
    % new data set  $D'$ .
  
```

Output: $H(x) = h'(h_1(x), \dots, h_T(x))$

多响应线性回归(MLR)作为次级学习器的学习算法 效果较好

贝叶斯模型平均(BMA)

