

机器学习

04 集成学习

李祎

liyi@dlut.edu.cn



大连理工大学 人工智能学院

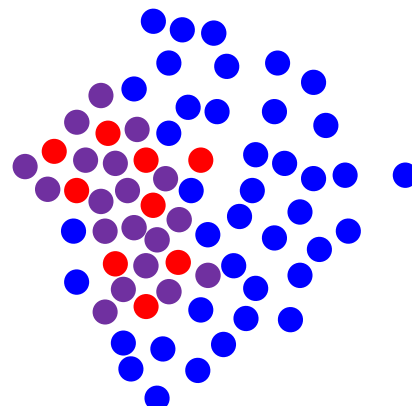
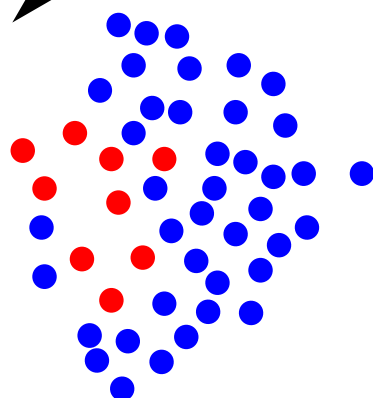
School of Artificial Intelligence, Dalian University of Technology

- 什么是集成学习
- Bagging与随机森林
- Boosting: Adaboost
- 结合策略
- 多样性

类别不平衡问题



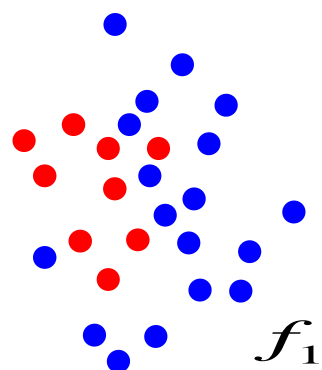
➤ 过采样 (oversampling):



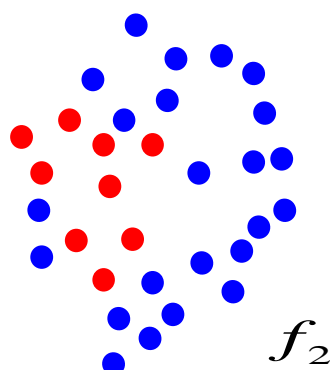
- ✓ 样本复制
- ✓ 样本插值
- ✓ 样本生成 (GAN)

[1] Chawla N V, Bowyer K W, et al. **SMOTE: Synthetic Minority Over-Sampling Technique**. *JAIR*, 2002.

➤ 欠采样 (undersampling)

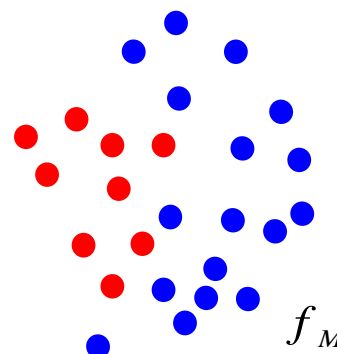


f_1



f_2

...



f_M

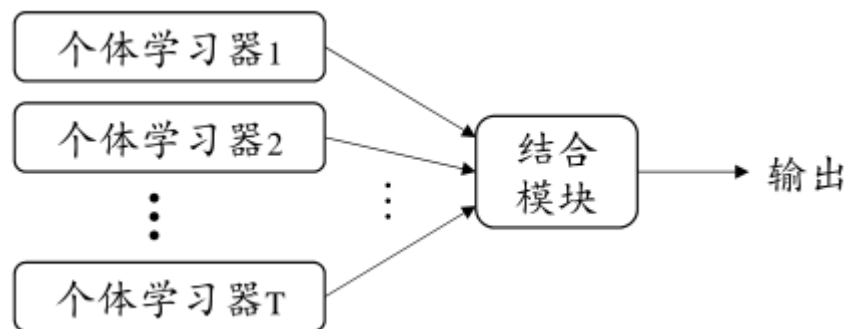
- 集成学习

$$f = \frac{1}{M} \sum_{m=1}^M f_m$$

- ✓ EasyEnsemble^[2]
- ✓ BalanceCascade^[2]

[2] Xu-Ying Liu, Jianxin Wu, Zhi-Hua Zhou. **Exploratory Undersampling for Class-Imbalance Learning**. *IEEE TSMCB*, 2009.

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



都是同一类型学习器：同质 基学习器 基学习算法
包含不同类型学习器：异质 组件学习器

□ 例如：一共有**500**个样本数据，每个子模型只用**100**个样本数据进行训练。假如每个子模型只有**51 %**的准确率（至少要比随机猜一个**50%**要高一点），

□ 如果只有一个子模型，整体准确率：**51%**

□ 如果有三个子模型，整体准确率： $0.51^3 + C_3^2 \cdot 0.51^2 \cdot 0.49 = 51.5\%$

□ 如果有**500**个子模型，整体准确率： $\sum_{i=251}^{500} C_{500}^i \cdot 0.51^i \cdot 0.49^{500-i} = 65.6\%$

□ 如果子模型的准确率**60%**，**500**个子模型集成：

$$\sum_{i=251}^{500} C_{500}^i \cdot 0.6^i \cdot 0.4^{500-i} = 99.999\%$$

- 如何获得比单一学习器更好的性能？
- 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_1	√	×	×
h_2	×	√	√	h_2	√	√	×	h_2	×	√	×
h_3	√	×	√	h_3	√	√	×	h_3	×	×	√
集群	√	√	√	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			

- 集成个体应该：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

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

- 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
- 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
 - 个体学习器间**存在**强依赖关系、必须**串行**生成的序列化方法：
Boosting
 - 个体学习器间**不存在**强依赖关系、可同时生成的**并行**化方法：
Bagging

Bagging与随机森林

- 并行式集成学习方法最著名的代表
- 基本思想：并行构建 T 个分类器、并行训练，最终将所有分类器的结果进行综合（平均或投票），达到最终的预测结果。
- Bagging: Bootstrap AGGregatING
- 自主采样法 (bootstrap sampling)

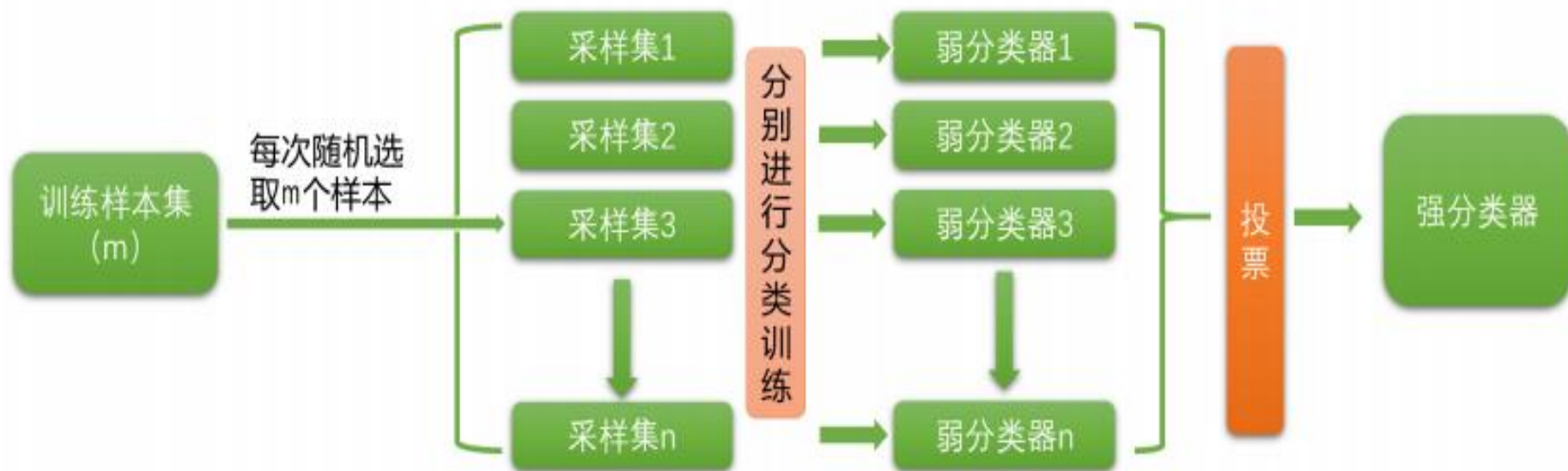
- 给定包含 m 个样本的数据集 D , 采样产生数据集 D' :
- 每次随机从 D 中挑选一个样本, 将其拷贝放入 D' , 然后再将该样本放回初始数据集 D 中, 使得该样本在下次采样时仍有可能被采到;
- 这个过程重复执行 m 次后, 我们就得到了包含 m 个样本的数据集 D' 。
- D 中有一部分样本会在 D' 中多次出现, 而另一部分样本不出现。可以做一个简单的估计, 样本在 m 次采样中始终不被采到的概率是 $(1 - \frac{1}{m})^m$, 取极限得到:

$$\lim_{m \rightarrow \infty} (1 - \frac{1}{m})^m \rightarrow \frac{1}{e} \approx 0.368$$

Bagging



- 采样出 T 个含 m 个训练样本的采样集
- 基于每个采样集训练出一个基学习器
- 再将这些基学习器进行结合



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

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

□ 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

□ 直接适用于多分类、回归等任务

□ 可使用包外估计

- 由于基学习器只使用了初始训练集中约63.2%的样本，剩下的约36.8%的样本可用作验证集来对泛化性能进行“包外估计”（out-of-bag estimate）。

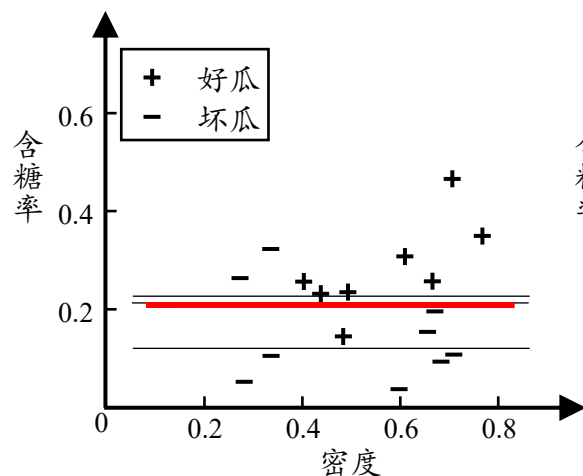
- 需要记录每个基学习器使用的训练样本 D_t
- $H^{oob}(x)$ 表示对样本 x 的包外预测，即仅考虑那些未使用样本 x 训练的基学习器在 x 上的预测

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

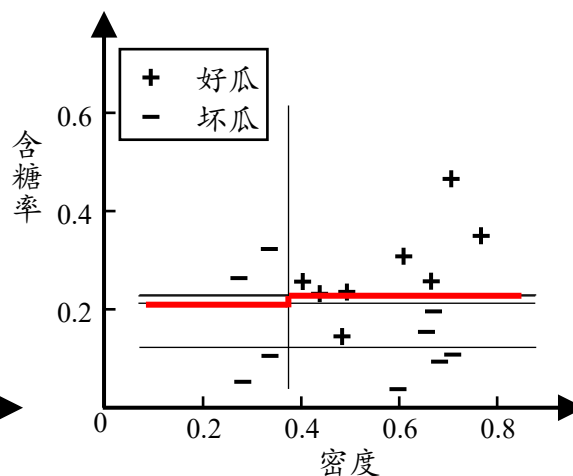
- Bagging泛化误差的包外估计为：

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

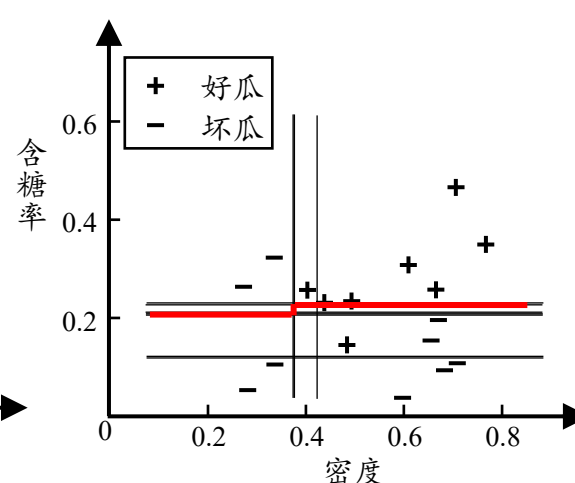
Bagging实验结果



(a) 3个基学习器



(b) 5个基学习器



(c) 11个基学习器

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

偏差-方差分解



对回归任务，泛化误差可通过“偏差-方差分解”拆解为：

$$E(f; D) = \underbrace{bias^2(x)}_{\text{red}} + \underbrace{var(x)}_{\text{blue}} + \underbrace{\varepsilon^2}_{\text{green}}$$

期望输出与真实
输出的差别

$$bias^2(x) = (\bar{f}(x) - y)^2$$

同样大小的训练集
的变动，所导致的
性能变化

$$var(x) = \mathbb{E}_D \left[(f(x; D) - \bar{f}(x))^2 \right]$$

训练样本的标记与
真实标记有区别

表达了当前任务上任何学习算法
所能达到的期望泛化误差下界

$$\varepsilon^2 = \mathbb{E}_D \left[(y_D - y)^2 \right]$$

泛化性能是由学习算法的能力、数据的充分性以及学习任务本身的难度共同决定

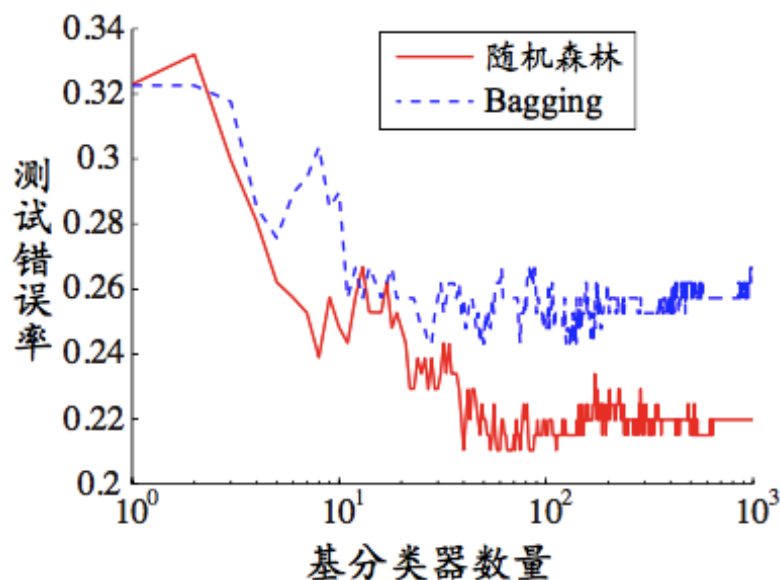
- 随机森林(Random Forest, 简称RF)是bagging的一个扩展变种
- 采样的随机性
- 属性选择的随机性
 - 传统决策树：

当前结点的所有属性中选择一个最优属性。
 - RF：

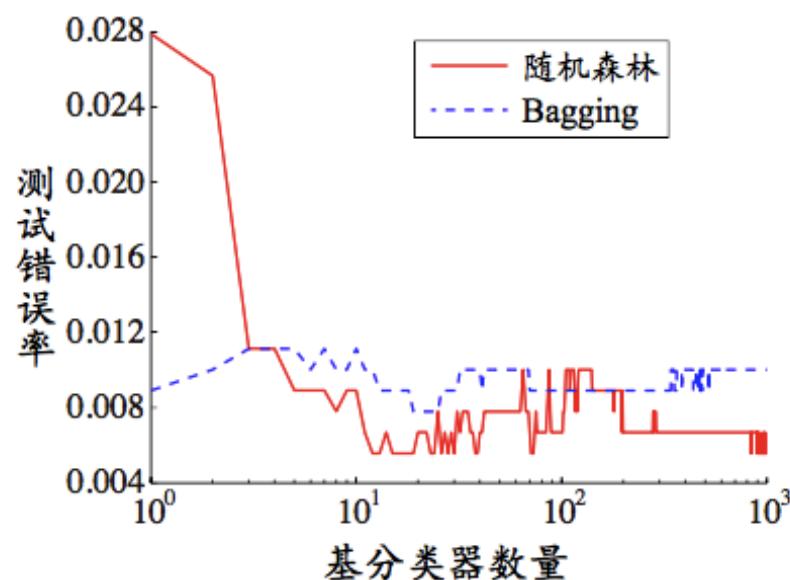
对基决策树的每个结点，
先从该结点的属性集合中随机选择一个包含 k 个属性的子集，
然后再从这个子集中选择一个最优属性用于划分。

基学习器的多样性:

- Bagging 中仅通过**样本扰动**产生,
- 随机森林中增加了**属性扰动**, 使得最终性能通过个体学习器之间差异度进一步提升。



(a) glass 数据集



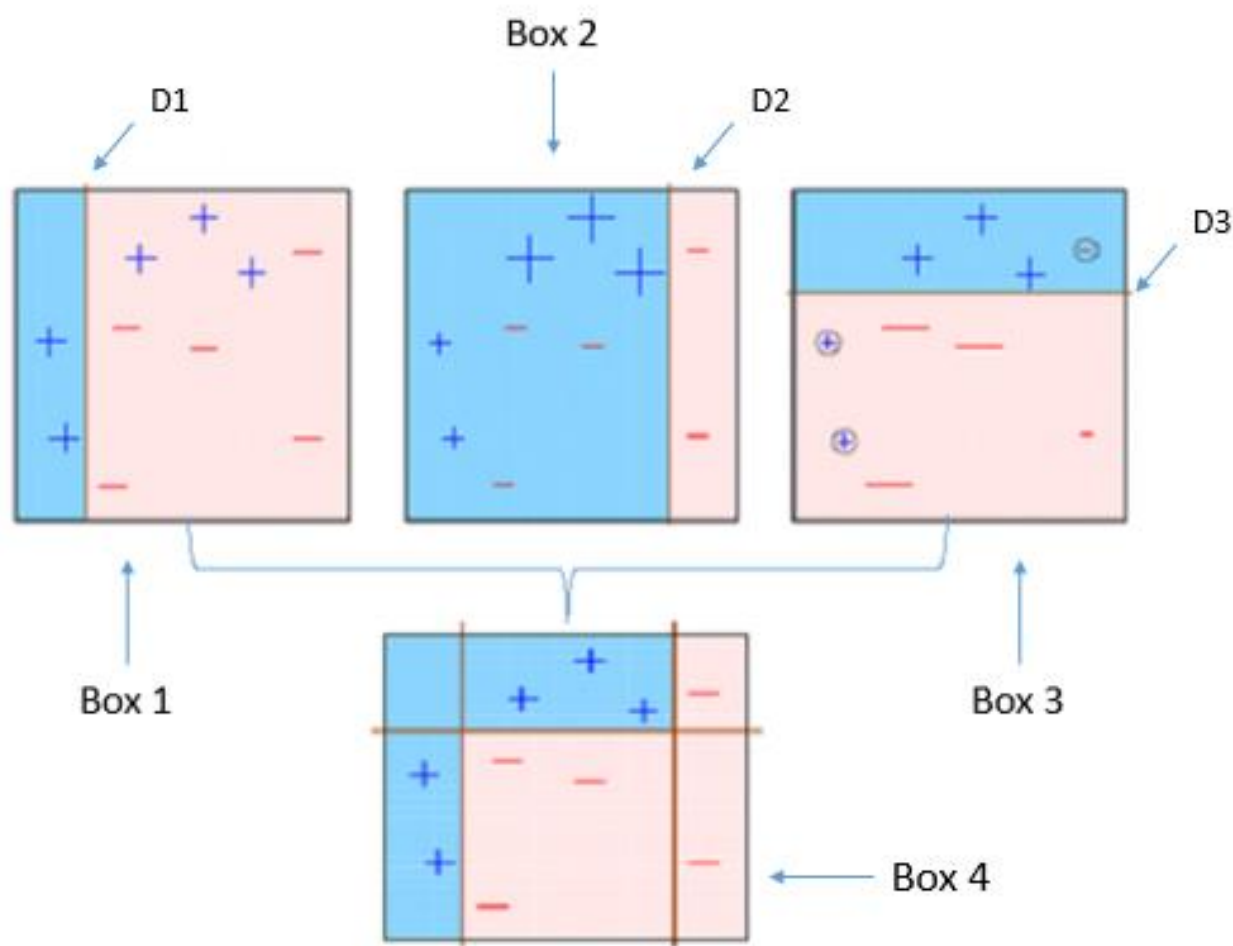
(b) auto-mpg 数据集

Boosting

Boosting



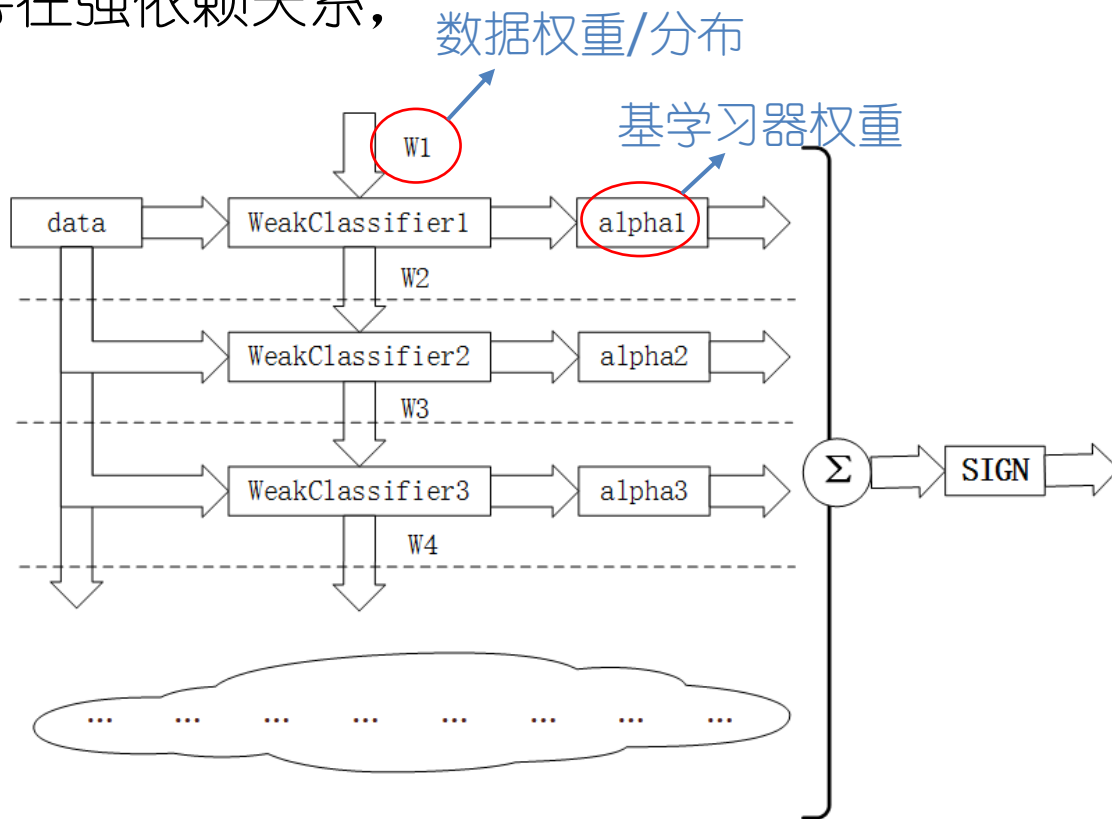
- 个体学习器存在强依赖关系，
- 串行生成
- 每次调整训练数据的样本分布



Boosting



- 个体学习器存在强依赖关系，
- 串行生成



1. 给定初始训练数据，由此训练出第一个基学习器，并确定其权重；
2. 对样本分布进行调整，在之前学习器做错的样本上投入更多关注；
3. 用调整后的样本，训练下一个基学习器，确定其权重；
4. 重复上述过程 T 次，将 T 个学习器加权结合。

□ 给定数据集 $D = \{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_m, y_m)\}$

其中, $y_i \in \{-1, +1\}$

□ 模型输出: 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: 0/1损失的一致性替代损失函数

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

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

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

$$H(\mathbf{x}) = \frac{1}{2} \ln \frac{P(f(\mathbf{x}) = 1 \mid \mathbf{x})}{P(f(\mathbf{x}) = -1 \mid \mathbf{x})}$$

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

指数损失函数最小化，则分类错误率也将最小化，说明指数损失函数是分类任务原来0/1损失函数的一致替代函数。

- 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

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

- 令指数损失函数的导数为0，即

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

- 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

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

□ 于是，理想的基学习器：

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

□ 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：

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

□ 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 $f(x), h(x) \in \{-1, +1\}$ 有：

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

□ 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

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

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: $\mathcal{D}_1(\mathbf{x}) = 1/m$. 初始化样本权重

2: **for** $t = 1, 2, \dots, T$ **do**

3: $h_t = \mathcal{L}(D, \mathcal{D}_t)$;

4: $\epsilon_t = P_{\mathbf{x} \sim \mathcal{D}_t}(h_t(\mathbf{x}) \neq f(\mathbf{x}))$; h_t 错误率

5: **if** $\epsilon_t > 0.5$ **then break**

6: $\alpha_t = \frac{1}{2} \ln \left(\frac{1-\epsilon_t}{\epsilon_t} \right)$; 确定 h_t 权重

7: $\mathcal{D}_{t+1}(\mathbf{x}) = \frac{\mathcal{D}_t(\mathbf{x})}{Z_t} \times \begin{cases} \exp(-\alpha_t), & \text{if } h_t(\mathbf{x}) = f(\mathbf{x}) \\ \exp(\alpha_t), & \text{if } h_t(\mathbf{x}) \neq f(\mathbf{x}) \end{cases}$ 更新样本分布
$$= \frac{\mathcal{D}_t(\mathbf{x}) \exp(-\alpha_t f(\mathbf{x}) h_t(\mathbf{x}))}{Z_t}$$

8: **end for**

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

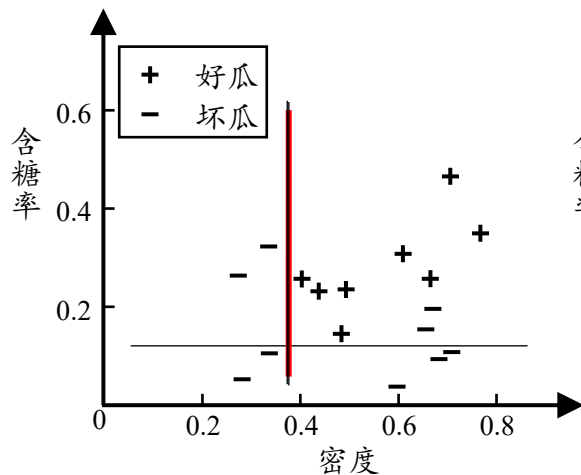
□ 优点

- 不容易发生过拟合；
- 由于AdaBoost并没有限制弱学习器的种类，所以可以使用不同的学习算法来构建弱分类器；
- 相对于bagging算法和Random Forest算法，AdaBoost充分考虑的每个分类器的权重；
- AdaBoost的参数少，实际应用中不需要调节太多的参数。

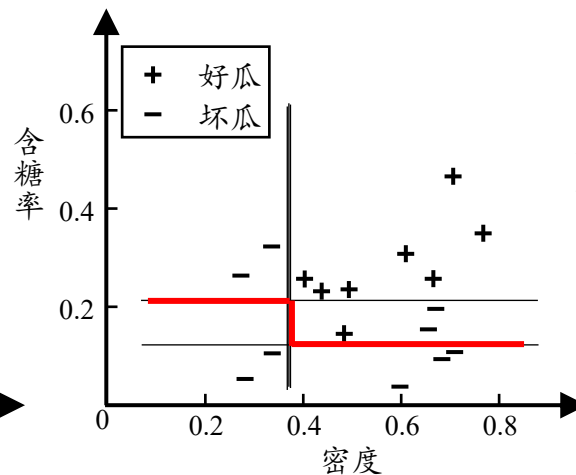
□ 缺点

- AdaBoost迭代次数（弱分类器数目）不好设定，可以使用交叉验证来确定；
- 对异常样本敏感，异常样本在迭代中可能会获得较高的权重，影响最终的强学习器的预测准确性；
- 训练比较耗时。

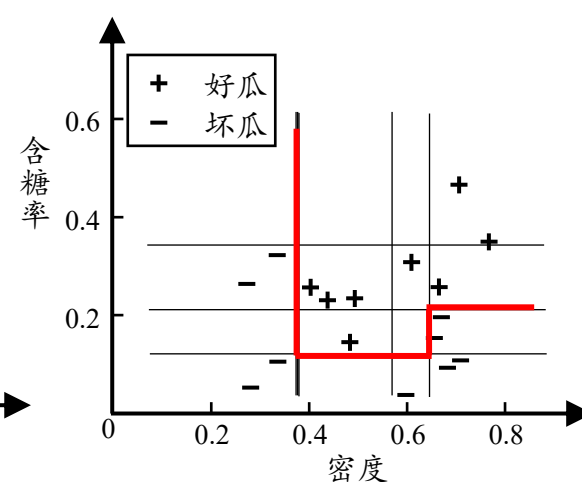
AdaBoost实验



(a) 3个基学习器



(b) 5个基学习器



(c) 11个基学习器

□ 从偏差-方差的角度：boosting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

结合策略

□ 简单平均法

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

□ 加权平均法

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

- 简单平均法是加权平均法的特例
- 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
 - 不同方式确定基学习器权重
- 加权平均法未必一定优于简单平均法
 - 训练样本不充分或含有噪声，导致学出的权重不完全可靠
- 一般而言，个体学习器性能相差较大：加权平均
性能相近：简单平均法

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)

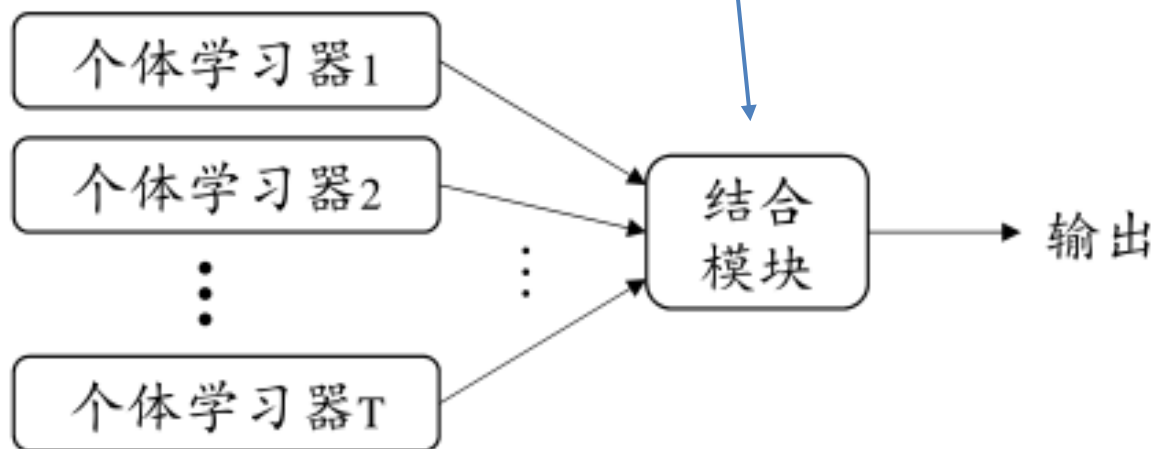
$$H(\mathbf{x}) = c_{\arg \max_j \sum_{i=1}^T h_i^j(\mathbf{x})}$$

不同类型的 $h_i^j(x)$
值不能混用

加权投票法 (weighted voting)

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

- 初级学习器+次级学习器/元学习器
- Stacking是学习法的典型代表



多样性

- 定义学习器 h_i 的分歧(ambiguity):

$$A(h_i | \mathbf{x}) = (h_i(\mathbf{x}) - H(\mathbf{x}))^2$$

- 集成的分歧:

$$\begin{aligned}\bar{A}(h | \mathbf{x}) &= \sum_{i=1}^T w_i A(h_i | \mathbf{x}) \\ &= \sum_{i=1}^T w_i (h_i(\mathbf{x}) - H(\mathbf{x}))^2\end{aligned}$$

多样性：误差-分歧分解



- 分歧项代表了个体学习器在样本 x 上的不一致性，即在一定程度上反映了个体学习器的多样性，个体学习器 h_i 和集成 H 的平方误差分别为

$$E(h_i | \mathbf{x}) = (f(\mathbf{x}) - h_i(\mathbf{x}))^2$$

$$E(H | \mathbf{x}) = (f(\mathbf{x}) - H(\mathbf{x}))^2$$

多样性：误差-分歧分解



□ 令 $\bar{E}(h | \mathbf{x}) = \sum_{i=1}^T w_i \cdot E(h_i | \mathbf{x})$ 表示个体学习器误差的加权均值，有

$$\begin{aligned}\bar{A}(h | \mathbf{x}) &= \sum_{i=1}^T w_i E(h_i | \mathbf{x}) - E(H | \mathbf{x}) \\ &= \bar{E}(h | \mathbf{x}) - E(H | \mathbf{x}) .\end{aligned}$$

□ 上式对所有样本 \mathbf{x} 均成立，令 $p(\mathbf{x})$ 表示样本的概率密度，则在全样本上有

$$\sum_{i=1}^T w_i \int A(h_i | \mathbf{x}) p(\mathbf{x}) d\mathbf{x} = \sum_{i=1}^T w_i \int E(h_i | \mathbf{x}) p(\mathbf{x}) d\mathbf{x} - \int E(H | \mathbf{x}) p(\mathbf{x}) d\mathbf{x}$$

多样性：误差-分歧分解



- 个体学习器 h_i 在全样本上的泛化误差和分歧项分别为：

$$E_i = \int E(h_i | \mathbf{x})p(\mathbf{x})d\mathbf{x}$$

$$A_i = \int A(h_i | \mathbf{x})p(\mathbf{x})d\mathbf{x}$$

- 集成的泛化误差为：

$$E = \int E(H | \mathbf{x})p(\mathbf{x})d\mathbf{x}$$

- 令 $\overline{E} = \sum_{i=1}^T w_i E_i$ 表示个体学习器泛化误差的加权均值，
 $\overline{A} = \sum_{i=1}^T w_i A_i$ 表示个体学习器的加权分歧值，有

$$E = \overline{E} - \overline{A}$$

多样性：误差-分歧分解



- 令 $\bar{E} = \sum_{i=1}^T w_i E_i$ 表示个体学习器泛化误差的加权均值，
 $\bar{A} = \sum_{i=1}^T w_i A_i$ 表示个体学习器的加权分歧值，有

$$E = \bar{E} - \bar{A}$$

- 这个漂亮的式子显示：个体学习器精确性越高、多样性越大，则集成效果越好。称为误差-分歧分解
- 为什么不能直接把 $\bar{E} - \bar{A}$ 作为优化目标来求解？
- 现实任务中很难直接对 $\bar{E} - \bar{A}$ 进行优化，
 - 它们定义在整个样本空间上
 - \bar{A} 不是一个可直接操作的多样性度量
 - 上面的推导过程只适用于回归学习，难以直接推广到分类学习任务上去

- 多样性度量(diversity measure)用于度量集成中个体学习器的多样性
- 对于二分类问题，分类器 h_i 与 h_j 的预测结果联立表(contingency table)为

	$h_i = +1$	$h_i = -1$
$h_j = +1$	a	c
$h_j = -1$	b	d

$$a + b + c + d = m$$

□ 常见的多样性度量

- 不合度量(Disagreement Measure)

$$dis_{ij} = \frac{b + c}{m}$$

- 相关系数(Correlation Coefficient)

$$\rho_{ij} = \frac{ad - bc}{\sqrt{(a + b)(a + c)(c + d)(b + d)}}$$

□ 常见的多样性度量

- Q-统计量 (Q-Statistic)

$$Q_{ij} = \frac{ad - bc}{ad + bc} \quad |Q_{ij}| \leq |\rho_{ij}|$$

- K-统计量 (Kappa-Statistic)

$$\kappa = \frac{p_1 - p_2}{1 - p_2} \quad \begin{aligned} p_1 &= \frac{a + d}{m}, \\ p_2 &= \frac{(a + b)(a + c) + (c + d)(b + d)}{m^2} \end{aligned}$$

□ 常见的增强个体学习器的多样性的方法

- 数据样本扰动
- 输入属性扰动
- 输出表示扰动
- 算法参数扰动

□ 不同的多样性增强机制可以同时使用

- 例如，随机森林

□ 数据样本扰动通常是基于采样法

- Bagging中的自助采样法
- Adaboost中的序列采样

数据样本扰动对“不稳定基学习器”很有效

□ 对数据样本的扰动敏感的基学习器(不稳定基学习器)

- 决策树，神经网络等

□ 对数据样本的扰动不敏感的基学习器(稳定基学习器)

- 线性学习器，支持向量机，朴素贝叶斯，k近邻等

- 什么是集成学习
- Bagging与随机森林
- Boosting
 - Adaboost
- 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