



第8章 集成学习



李政伟

-----• 中国矿业大学 计算机科学与技术学院 •-----



8.1 个体与集成

- **集成学习**(ensemble learning): 通过构建并结合 **多个学习器** 来完成学习任务, 有时候也被称为 **多分类器系统**(multi-classifier system)、**基于委员会的学习**(committee voting method)。
- **一般结构**: 先产生一组“个体学习器”(individual learner), 再用某种策略将它们结合起来。
- **个体学习器**: 通常由现有的学习算法从训练数据产生, 如C4.5决策树算法、BP神经网络算法等。

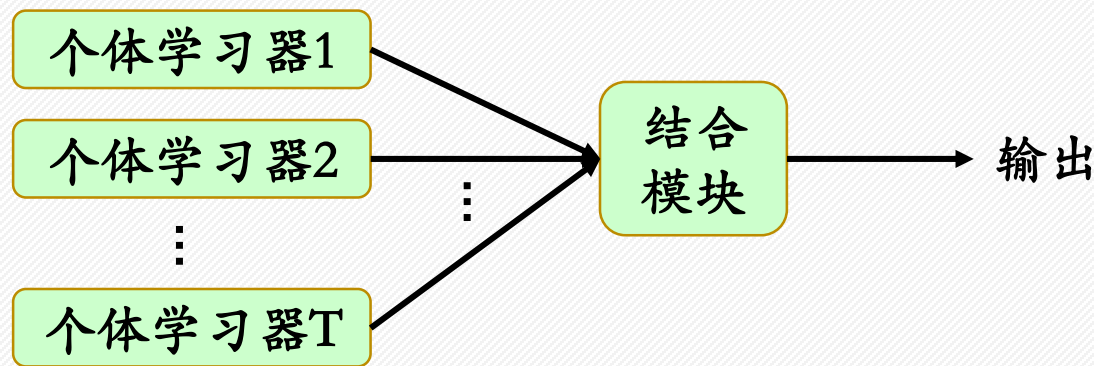


图8.1 集成学习示意图



8.1 个体与集成

同质(homogeneous)集成

- 集成中只包含**同种**类型的个体学习器。
- 个体学习器：又称为**基学习器**(base learner)。
- 对应学习算法：基学习算法(base learning algorithm)。

异质 (heterogeneous)集成

- 集成中包含**不同类型**的个体学习器。
- 个体学习器：由不同的学习算法生成，不再有基学习算法。
- 个体学习器：常称为**组件学习器**(component learner)。

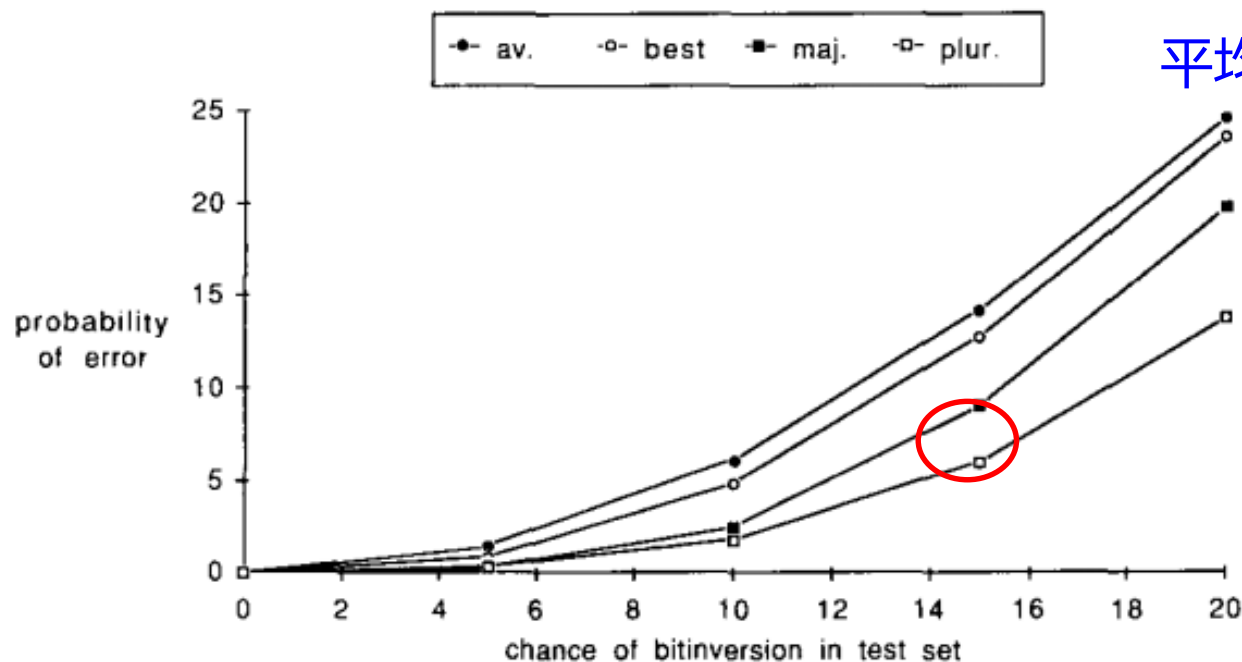


Why Ensemble?

集成的泛化性能通常显著优于单个学习器的泛化性能

一个观察：

误差 (曲线越低越好)



平均神经网络的错误率

最优神经网络的错误率

两种简单的神经网络集成算法的错误率

Fig. 4. Performance versus noise level in the test set is shown for individual and for consensus decisions. Data displayed shows the average and the best network, as well as collective decisions using majority and plurality for seven networks trained on individual training sets.

[Hansen & Salamon, TPAMI90]



如何得到好的集成？

	测例1	测例2	测例3
h_1	√	√	×
h_2	×	√	√
h_3	√	×	√
集成	√	√	√

(a)集成提升性能

	测例1	测例2	测例3
h_1	√	√	×
h_2	√	√	×
h_3	√	√	×
集成	√	√	×

(b)集成不起作用

	测例1	测例2	测例3
h_1	√	×	×
h_2	×	√	×
h_3	×	×	√
集成	×	×	×

(c)集成起负作用

集成的个体学习器应 “**好而不同**”

现实各类机器学习、数据挖掘应用中，广泛使用集成学习技术
想获胜，用集成！



8.1 个体与集成

考虑二分类问题 $y \in \{-1, +1\}$ 和真实函数 f ，假定基分类器的错误率为 ϵ ，对每个基分类器 h_i

$$P(h_i(\mathbf{x}) \neq f(\mathbf{x})) = \epsilon \quad (8.1)$$

假设集成通过简单投票法结合 T 个基分类器，若有超过半数的基分类器正确，则集成分类就正确：

$$H(\mathbf{x}) = \text{sign}\left(\sum_{i=1}^T h_i(\mathbf{x})\right) \quad (8.2)$$

假设基分类器的错误率相互独立，则集成的错误率为

$$\begin{aligned} P(h_i(\mathbf{x}) \neq f(\mathbf{x})) &= \sum_{k=0}^{\lfloor T/2 \rfloor} \binom{T}{k} (1 - \epsilon)^k \epsilon^{T-k} \\ &\leq \exp\left(-\frac{1}{2} T (1 - 2\epsilon)^2\right) \end{aligned} \quad (8.3)$$

很多成功的集成学习方法

■ 序列化方法：个体学习器间存在强依赖关系

- **AdaBoost** [Freund & Schapire, JCSS97]
- GradientBoost [Friedman, AnnStat01]
- LPBoost [Demiriz, Bennett, Shawe-Taylor, MLJ06]
-

■ 并行化方法：个体学习器间不存在强依赖关系

- **Bagging** [Breiman, MLJ96]
- **Random Forest** [Breiman, MLJ01]
- Random Subspace [Ho, TPAMI98]
-



8.2 Boosting

Boosting: 一族可将弱学习器提升为强学习器的算法。

工作机制: 先从初始训练集训练出一个基学习器, 再根据基学习器的表现对训练样本分布进行调整, 使得先前基学习器分错的训练样本在后续受到更多关注, 然后基于调整后的样本分布来训练下一个基学习器; 如此重复进行, 直至基学习器达到事先指定的 T 个, 最终将这 T 个基学习器进行加权结合。

$$h_1(x) \in \{-1, 1\}$$

$$h_2(x) \in \{-1, 1\}$$

$$\vdots$$

$$h_T(x) \in \{-1, 1\}$$

Weak classifiers

$$H_T(x) = \text{sign}\left(\sum_{t=1}^T \alpha_t h_t(x)\right)$$

Strong classifier



Adaboost的提出

- 1990年，Schapire最先构造出一种多项式级的算法，即最初的Boost算法；
- 1993年，Drunker和Schapire第一次将神经网络作为弱学习器，应用Boosting算法解决OCR问题；
- 1995年，Freund和Schapire提出了Adaboost(Adaptive Boosting)算法，效率和原来Boosting算法一样，但是不需要任何关于弱学习器性能的先验知识，可以非常容易地应用到实际问题中。

AdaBoost

Adaptive Boosting

A learning algorithm

Building a **strong** classifier based on a lot of **weaker** ones



Adaboost基本概念

- **关键：** 如何解决好两个问题。
 1. **每一轮如何改变训练数据的权值或概率分布？**
 - **措施：** 提高那些被前一轮弱分类器错误分类样本的权值，降低那些被正确分类样本的权值
 2. **如何将弱分类器组合成一个强分类器？**
 - **措施：** 加权多数表决，加大分类误差率小的弱分类器的权值，使其在表决中起较大的作用，减小分类误差率大的弱分类器的权值，使其在表决中起较小的作用。



AdaBoost算法

- **输入**：训练数据集 $D = \{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_m, y_m)\}$, $\mathbf{x}_i \in \mathcal{X} \subseteq \mathbb{R}^d$, $y_i \in \mathcal{Y} = \{-1, +1\}$

- **输出**：最终分类器 $H(\mathbf{x})$

(1) 设定弱分类器数目为 T , 令 $t = 1$, 使用**均匀分布**初始化训练样本集的**权重分布**, 令 m 维向量 \mathcal{D}_1 表示第1次需要更新的**样本权重**, \mathcal{D}_1 的初始值设为

$$\mathcal{D}_t = (\mathcal{D}_{t1}, \mathcal{D}_{t2}, \dots, \mathcal{D}_{tm})^T = \left(\frac{1}{m}, \frac{1}{m}, \dots, \frac{1}{m}\right)^T$$

(2) 使用**权重分布**为 \mathcal{D}_t 的训练样本集 D 学习得到第 t 个弱学习器 h_t

(3) 计算 h_t 在训练样本集 D 上的分类错误率 ϵ_t

$$\epsilon_t = \sum_{i=1}^m \mathcal{D}_{ti} \mathbb{I}(h_t(\mathbf{x}_i) \neq y_i)$$

如果 $\epsilon_t > 0.5$, 则舍弃。



AdaBoost算法

(4) 确定弱学习器 h_t 的组合权重 α_t 。由于弱学习器 h_t 的权重与其分类性能相关，分类错误率 ϵ_t 越小的 h_t ，则其权重 α_t 应该越大，有

$$\alpha_t = \frac{1}{2} \ln \frac{1 - \epsilon_t}{\epsilon_t}$$

(5) 依据弱学习器 h_t 对训练样本集 D 的分类错误率 ϵ_t ，更新样本权重

$$\mathcal{D}_{t+1,i} = \frac{\mathcal{D}_{t,i} \exp(-\alpha_t y_i h_t(x_i))}{Z_t}$$

其中， $Z_t = \sum_{i=1}^N \mathcal{D}_{t,i} \exp(-\alpha_t y_i h_t(x_i))$ 为归一化因子，保证更新后权重向量为概率分布。

(6) 若 $t < T$ ，则令 $t = t + 1$ ，并返回步骤(2)，否则，执行步骤(7)；

(7) 对于 T 个分类器 h_1, h_2, \dots, h_T ，分别将每个按权重 α_t 进行组合，得最终分类器。

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



AdaBoost算法推导

AdaBoost——加性模型(additive model)

基学习器的线性组合

$$H(x) = \sum_{t=1}^T \alpha_t h_t(x) \quad (8.4)$$

最小化指数损失函数

$$\ell_{exp}(H|\mathcal{D}) = E_{x \sim \mathcal{D}}[e^{-f(x)H(x)}] \quad (8.5)$$

若 $H(x)$ 能令指数损失函数最小化，则式(8.5)对 $H(x)$ 的偏导

$$\frac{\partial \ell_{exp}}{\partial H(x)} = -e^{-H(x)} P(f(x) = 1|x) + e^{H(x)} P(f(x) = -1|x) \quad (8.6)$$

令式(8.6)为零可解得

$$H(x) = \frac{1}{2} \ln \frac{P(f(x) = 1|x)}{P(f(x) = -1|x)} \quad (8.7)$$



AdaBoost算法说明

因此,

$$\begin{aligned}\text{sign}(H(\boldsymbol{x})) &= \text{sign}\left(\frac{1}{2} \ln \frac{P(f(\boldsymbol{x}) = 1|\boldsymbol{x})}{P(f(\boldsymbol{x}) = -1|\boldsymbol{x})}\right) \\ &= \begin{cases} 1, P(f(\boldsymbol{x}) = 1|\boldsymbol{x}) > P(f(\boldsymbol{x}) = -1|\boldsymbol{x}) \\ -1, P(f(\boldsymbol{x}) = 1|\boldsymbol{x}) < P(f(\boldsymbol{x}) = -1|\boldsymbol{x}) \end{cases} \\ &= \arg \max_{y \in \{-1, 1\}} P(f(\boldsymbol{x}) = y|\boldsymbol{x}) \end{aligned} \tag{8.8}$$

这意味着 $\text{sign}(H(\boldsymbol{x}))$ 达到了**贝叶斯最优错误率**。换言之, 若指数损失函数最小化, 则**分类错误率也将最小化**。



AdaBoost算法说明

$$\begin{aligned}\ell_{exp}(\alpha_t h_t | \mathcal{D}_t) &= E_{x \sim \mathcal{D}_t} [e^{-f(x) \alpha_t h_t(x)}] \\ &= E_{x \sim \mathcal{D}_t} [e^{-\alpha_t} \mathbb{I}(f(x) = h_t(x)) + e^{\alpha_t} \mathbb{I}(f(x) \neq h_t(x))] \\ &= e^{-\alpha_t} P_{x \sim \mathcal{D}_t}(f(x) = h_t(x)) + e^{\alpha_t} P_{x \sim \mathcal{D}_t}(f(x) \neq h_t(x)) \\ &= e^{-\alpha_t} (1 - \epsilon_t) + e^{\alpha_t} \epsilon_t\end{aligned} \tag{8.9}$$

其中 $\epsilon_t = P_{x \sim \mathcal{D}_t}(f(x) \neq h_t(x))$ 考虑指数损失函数的导数

$$\frac{\partial \ell_{exp}(\alpha_t h_t | \mathcal{D}_t)}{\partial \alpha_t} = -e^{-\alpha_t} (1 - \epsilon_t) + e^{\alpha_t} \epsilon_t \tag{8.10}$$

令式(8.10)为零可解得

$$\alpha_t = \frac{1}{2} \ln \frac{1 - \epsilon_t}{\epsilon_t} \tag{8.11}$$



AdaBoost算法说明

$$\begin{aligned}\ell_{exp}(H_{t-1} + h_t | \mathcal{D}) &= \mathbb{E}_{x \sim \mathcal{D}}[e^{-f(x)(H_{t-1}(x) + h_t(x))}] \\ &= \mathbb{E}_{x \sim \mathcal{D}}[e^{-f(x)H_{t-1}(x)} e^{-f(x)h_t(x)}]\end{aligned}\tag{8.12}$$

注意到 $f^2(x) = h_t^2(x) = 1$ 式(8.12)可使用 $e^{-f(x)h_t(x)}$ 的泰勒展式近似为

$$\begin{aligned}\ell_{exp}(H_{t-1} + h_t | \mathcal{D}) &\approx \mathbb{E}_{x \sim \mathcal{D}}[e^{-f(x)H_{t-1}(x)} (1 - f(x)h_t(x) + \frac{f^2(x)h_t^2(x)}{2})] \\ &= \mathbb{E}_{x \sim \mathcal{D}}[e^{-f(x)H_{t-1}(x)} (1 - f(x)h_t(x) + \frac{1}{2})]\end{aligned}\tag{8.13}$$



AdaBoost算法说明

$$h_t(\mathbf{x}) = \arg \min_h \ell_{\exp}(H_{t-1} + h|\mathcal{D})$$

$$= \arg \min_h \mathbb{E}_{\mathbf{x} \sim \mathcal{D}} [e^{-f(\mathbf{x})H_{t-1}(\mathbf{x})} (1 - f(\mathbf{x})h(\mathbf{x}) + \frac{1}{2})]$$

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

$$= \arg \max_h \mathbb{E}_{\mathbf{x} \sim \mathcal{D}} \left[\frac{e^{-f(\mathbf{x})H_{t-1}(\mathbf{x})}}{\mathbb{E}_{\mathbf{x} \sim \mathcal{D}} [e^{-f(\mathbf{x})H_{t-1}(\mathbf{x})}]} f(\mathbf{x})h(\mathbf{x}) \right] \quad (8.14)$$

因 $\mathbb{E}_{\mathbf{x} \sim \mathcal{D}} [e^{-f(\mathbf{x})H_{t-1}(\mathbf{x})}]$ 是常数, 令

$$\mathcal{D}_t(\mathbf{x}) = \frac{\mathcal{D}(\mathbf{x})e^{-f(\mathbf{x})H_{t-1}(\mathbf{x})}}{\mathbb{E}_{\mathbf{x} \sim \mathcal{D}} [e^{-f(\mathbf{x})H_{t-1}(\mathbf{x})}]} \quad (8.15)$$

则根据数学期望的定义, 等价于令

$$h_t(\mathbf{x}) = \arg \max_h \mathbb{E}_{\mathbf{x} \sim \mathcal{D}} \left[\frac{e^{-f(\mathbf{x})H_{t-1}(\mathbf{x})}}{\mathbb{E}_{\mathbf{x} \sim \mathcal{D}} [e^{-f(\mathbf{x})H_{t-1}(\mathbf{x})}]} f(\mathbf{x})h(\mathbf{x}) \right] \quad (8.16)$$

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



AdaBoost算法说明

由 $f(x), h(x) \in \{-1, +1\}$, 有 $f(x)h(x) = 1 - 2\mathbb{I}(f(x) \neq h(x))$ (8.17)

则理想的基学习器

$$h_t(x) = \arg \min_h \mathbb{E}_{x \sim \mathcal{D}_t} [\mathbb{I}(f(x) \neq h(x))] \quad (8.18)$$

考虑到 \mathcal{D}_t 和 \mathcal{D}_{t+1} 的关系, 有

$$\begin{aligned} \mathcal{D}_{t+1}(x) &= \frac{\mathcal{D}(x)e^{-f(x)H_t(x)}}{\mathbb{E}_{x \sim \mathcal{D}}[e^{-f(x)H_t(x)}]} \\ &= \frac{\mathcal{D}(x)e^{-f(x)H_{t-1}(x)}e^{-f(x)\alpha_t h_t(x)}}{\mathbb{E}_{x \sim \mathcal{D}}[e^{-f(x)H_t(x)}]} \\ &= \mathcal{D}_t(x)e^{-f(x)\alpha_t h_t(x)} \frac{\mathbb{E}_{x \sim \mathcal{D}}[e^{-f(x)H_{t-1}(x)}]}{\mathbb{E}_{x \sim \mathcal{D}}[e^{-f(x)H_t(x)}]} \end{aligned} \quad (8.19)$$



Adaboost示例

x	0	1	2	3	4	5	6	7	8	9
y	1	1	1	-1	-1	-1	1	1	1	-1

均匀初始化 $\mathcal{D}_1 = (\mathcal{D}_{1,1}, \mathcal{D}_{1,2}, \dots, \mathcal{D}_{1,10})$ $\mathcal{D}_{1,i} = 0.1, \quad i = 1, 2, \dots, 10$

对于迭代过程 $t=1$

a、在权值分布为 \mathcal{D}_1 的数据集上，阈值取**2.5**，分类误差率**最小**，**基分类器** h_1 ：

$$h_1(x) = \begin{cases} 1, & x < 2.5 \\ -1, & x \geq 2.5 \end{cases}$$

b、 $h_1(x)$ 的误差率： $\epsilon_1 = P(h_1(x_i) \neq y_i) = 0.1 + 0.1 + 0.1 = 0.3$

c、 $h_1(x)$ 的系数： $\alpha_1 = \frac{1}{2} \log \frac{1 - \epsilon_1}{\epsilon_1} = 0.4236$

$$f_1(x) = 0.4236 h_1(x)$$



Adaboost示例（续）

x	0	1	2	3	4	5	6	7	8	9
y	1	1	1	-1	-1	-1	1	1	1	-1
D_1	0.1	0.1	0.1	0.1	0.1	0.1	0.1	0.1	0.1	0.1

确定阈值的取值及误差率

- 当阈值取2.5时，误差率为**0.3**。即 $x < 2.5$ 时取 1, $x > 2.5$ 时取 -1, 则数据6、7、8分错，误差率为0.3（简单理解：10个里面3个错的）
- 当阈值取5.5时，误差率最低为**0.4**。即 $x < 5.5$ 时取1, $x > 5.5$ 时取 -1, 则数据3、4、5、6、7、8分错，错误率为0.6>0.5, 故反过来，令 $x > 5.5$ 取 1, $x < 5.5$ 时取 -1, 则数据0、1、2、9分错，误差率为0.4
- 当阈值取8.5时，误差率为**0.3**。即 $x < 8.5$ 时取1, $x > 8.5$ 时取 -1, 则数据3、4、5分错，错误率为0.3
- 因此，阈值取2.5 或8.5时，误差率等值，所以可以任选一个作为基本分类器。这里**选2.5为例**。

Adaboost示例（续）

d、更新训练数据的权值分布

$$(\mathcal{D}_{2,1}, \mathcal{D}_{2,2}, \dots, \mathcal{D}_{2,10})$$

x	0	1	2	3	4	5	6	7	8	9
y	1	1	1	-1	-1	-1	1	1	1	-1
$h_1(x)$	1	1	1	-1	-1	-1	-1	-1	-1	-1
分类结果	对	对	对	对	对	对	错	错	错	对

$$\mathcal{D}_{2,i} = \frac{\mathcal{D}_1(x)}{Z_1} \times \begin{cases} \exp(-\alpha_1) & \text{if } h_t(x) = f(x) \\ \exp(\alpha_1) & \text{if } h_t(x) \neq f(x) \end{cases}$$

$$\mathcal{D}_2 = (0.07143, 0.07143, 0.07143, 0.07143, 0.07143, 0.07143, 0.16666, 0.16666, 0.16666, 0.07143)$$

x	0	1	2	3	4	5	6	7	8	9
y	1	1	1	-1	-1	-1	1	1	1	-1
\mathcal{D}_1	0.1	0.1	0.1	0.1	0.1	0.1	0.1	0.1	0.1	0.1
\mathcal{D}_2	0.07143	0.07143	0.07143	0.07143	0.07143	0.07143	0.16666	0.16666	0.16666	0.07143



Adaboost示例（续）

$$Z_1 = \sum_{i=1}^n D_{1i} \exp(-y_i \alpha_1 h_1(x_i))$$

y	1	1	1	-1	-1	-1	1	1	1	-1
$h_1(x)$	1	1	1	-1	-1	-1	-1	-1	-1	-1
结果	对	对	对	对	对	对	错	错	错	对
\mathcal{D}_1	0.1	0.1	0.1	0.1	0.1	0.1	0.1	0.1	0.1	0.1

$$= \sum_{i=1}^3 0.1 \times \exp(-[0.4263 \times 1 \times 1]) + \sum_{i=4}^{4-6,10} 0.1 \times \exp(-[0.4263 \times (-1) \times (-1)])$$

$$+ \sum_{i=7}^9 0.1 \times \exp(-[0.4263 \times 1 \times (-1)])$$

$$\approx 0.3928 + 0.4582 + 0.0655 = 0.9165$$

$$\mathcal{D}_{2,i} = \frac{\mathcal{D}_{1,i}}{Z_1} \exp(-y_i \alpha_1 h_1(x_i)) = \begin{cases} \frac{0.1}{0.9165} \exp(-[0.4263 \times 1 \times 1]) \\ \frac{0.1}{0.9165} \exp(-[0.4263 \times (-1) \times (-1)]) \\ \frac{0.1}{0.9165} \exp(-[0.4263 \times 1 \times (-1)]) \end{cases} \approx \begin{cases} 0.07143 & i = 1,2,3 \\ 0.07143 & i = 4,5,6,10 \\ 0.16666 & i = 7,8,9 \end{cases}$$

\mathcal{D}_2	0.07143	0.07143	0.07143	0.07143	0.07143	0.07143	0.16666	0.16666	0.16666	0.07143
-----------------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------



Adaboost示例 (续)

对迭代过程 $t=2$

x	0	1	2	3	4	5	6	7	8	9
y	1	1	1	-1	-1	-1	1	1	1	-1

a、在权值分布 \mathcal{D}_2 上，阈值 $v=8.5$ 时，分类误差率最低

$$h_2(x) = \begin{cases} 1, & x < 8.5 \\ -1, & x \geq 8.5 \end{cases}$$

b、误差率

$$\epsilon_2 = 0.07143 + 0.07143 + 0.07143 = 0.21429$$

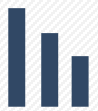
c、计算

$$\alpha_2 = \frac{1}{2} \log \frac{1 - \epsilon_2}{\epsilon_2} = 0.6496$$

$$f_2(x) = 0.6496 h_2(x)$$

d、更新权值分布

$$\mathcal{D}_3 = (0.04546, 0.04546, 0.04546, 0.16667, 0.16667, 0.16667, 0.10606, 0.10606, 0.10606, 0.04546)$$



Adaboost示例（续）

x	0	1	2	3	4	5	6	7	8	9
y	1	1	1	-1	-1	-1	1	1	1	-1
D_2	0.07143	0.07143	0.07143	0.07143	0.07143	0.07143	0.16666	0.16666	0.16666	0.07143

确定阈值的取值及误差率

- 当阈值取2.5时，误差率为0.49998。即 $x < 2.5$ 时取 1， $x > 2.5$ 时取 -1，则数据6、7、8分错，误差率为 $0.16666 * 3 = 0.49998$
- 当阈值取5.5时，误差率最低为0.28572。即 $x < 5.5$ 时取1， $x > 5.5$ 时取 -1，则数据3、4、5、6、7、8分错，错误率为 $0.07143 * 3 + 0.16666 * 3 = 0.71427 > 0.5$ ，故反过来，令 $x > 5.5$ 取 1， $x < 5.5$ 时取 -1，则数据0、1、2、9分错，误差率为 $0.07143 * 4 = 0.28572$
- 当阈值取8.5时，误差率为0.21429。即 $x < 8.5$ 时取1， $x > 8.5$ 时取 -1，则数据3、4、5分错，错误率为 $0.07143 * 3 = 0.21429$



Adaboost示例（续）

$$Z_2 = \sum_{i=1}^n D_{1i} \exp(-y_i \alpha_1 h_1(x_i))$$

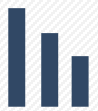
$$= \sum_{i=1}^3 0.07143 \times \exp(-[0.64963 \times 1 \times 1]) + \sum_{i=4}^{4-6} 0.07143 \times \exp(-[0.64963 \times (-1) \times 1])$$

$$+ \sum_{i=7}^9 0.16666 \times \exp(-[0.64963 \times 1 \times 1]) + \sum_{i=10}^{10-10} 0.07143 \times \exp(-[0.64963 \times (-1) \times (-1)])$$

$$\approx 0.11191 + 0.41033 + 0.26111 + 0.03730 = 0.82065$$

$$\mathcal{D}_{3,i} = \frac{D_{2,i}}{Z_2} \exp(-y_i \alpha_2 h_2(x_i)) = \begin{cases} \frac{0.07143}{0.82065} \exp(-[0.64963 \times 1 \times 1]) \\ \frac{0.07143}{0.82065} \exp(-[0.64963 \times (-1) \times 1]) \\ \frac{0.16666}{0.82065} \exp(-[0.64963 \times 1 \times 1]) \\ \frac{0.07143}{0.82065} \exp(-[0.64963 \times (-1) \times (-1)]) \end{cases} \approx \begin{cases} 0.04546 & i = 1,2,3 \\ 0.16667 & i = 4,5,6 \\ 0.10606 & i = 7,8,9 \\ 0.04546 & i = 10 \end{cases}$$

\mathcal{D}_2	0.07143	0.07143	0.07143	0.07143	0.07143	0.07143	0.16666	0.16666	0.16666	0.07143
\mathcal{D}_3	0.04546	0.04546	0.04546	0.16667	0.16667	0.16667	0.10606	0.10606	0.10606	0.04546



Adaboost示例（续）

x	0	1	2	3	4	5	6	7	8	9
y	1	1	1	-1	-1	-1	1	1	1	-1
\mathcal{D}_1	0.1	0.1	0.1	0.1	0.1	0.1	0.1	0.1	0.1	0.1
$h_1(x)$	1	1	1	-1	-1	-1	-1	-1	-1	-1
分类结果	对	对	对	对	对	对	错	错	错	对
\mathcal{D}_2	0.07143	0.07143	0.07143	0.07143	0.07143	0.07143	0.16666	0.16666	0.16666	0.07143
$h_2(x)$	1	1	1	1	1	1	1	1	1	-1
分类结果	对	对	对	错	错	错	对	对	对	对
\mathcal{D}_3	0.04546	0.04546	0.04546	0.16667	0.16667	0.16667	0.10606	0.10606	0.10606	0.04546



Adaboost示例（续）

对迭代过程 $t=3$

x	0	1	2	3	4	5	6	7	8	9
y	1	1	1	-1	-1	-1	1	1	1	-1

a、在权值分布 \mathcal{D}_3 上，阈值 $v=5.5$ 时，分类误差率最低

$$h_3(x) = \begin{cases} -1, & x < 5.5 \\ 1, & x \geq 5.5 \end{cases}$$

b、误差率

$$\epsilon_3 = 0.1820$$

c、计算

$$\alpha_3 = \frac{1}{2} \log \frac{1 - \epsilon_3}{\epsilon_3} = 0.7519$$

$$f_3(x) = 0.7519h_3(x)$$

d、更新权值分布

$$\mathcal{D}_4 = (0.1250, 0.1250, 0.1250, 0.1019, 0.1019, 0.1019, 0.0648, 0.0648, 0.0648, 0.1250)$$

最终分类器

$$H(x) = \text{sign}(f_3(x)) = \text{sign}[0.4236h_1(x) + 0.6496h_2(x) + 0.7519h_3(x)]$$

Adaboost示例 (续)

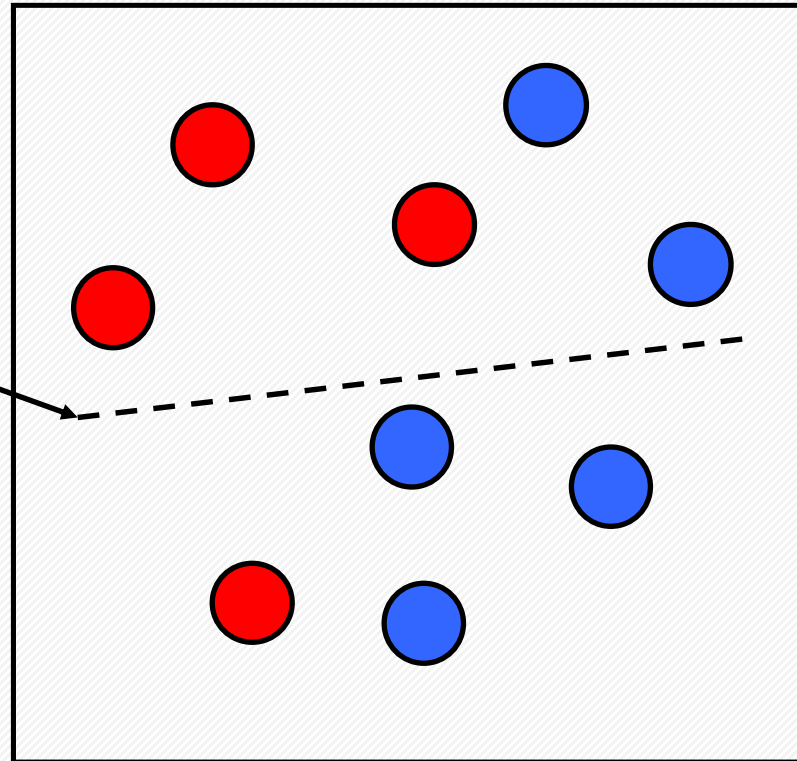
x	0	1	2	3	4	5	6	7	8	9
y	1	1	1	-1	-1	-1	1	1	1	-1
\mathcal{D}_3	0.04546	0.04546	0.04546	0.16667	0.16667	0.16667	0.10606	0.10606	0.10606	0.04546

确定阈值的取值及误差率

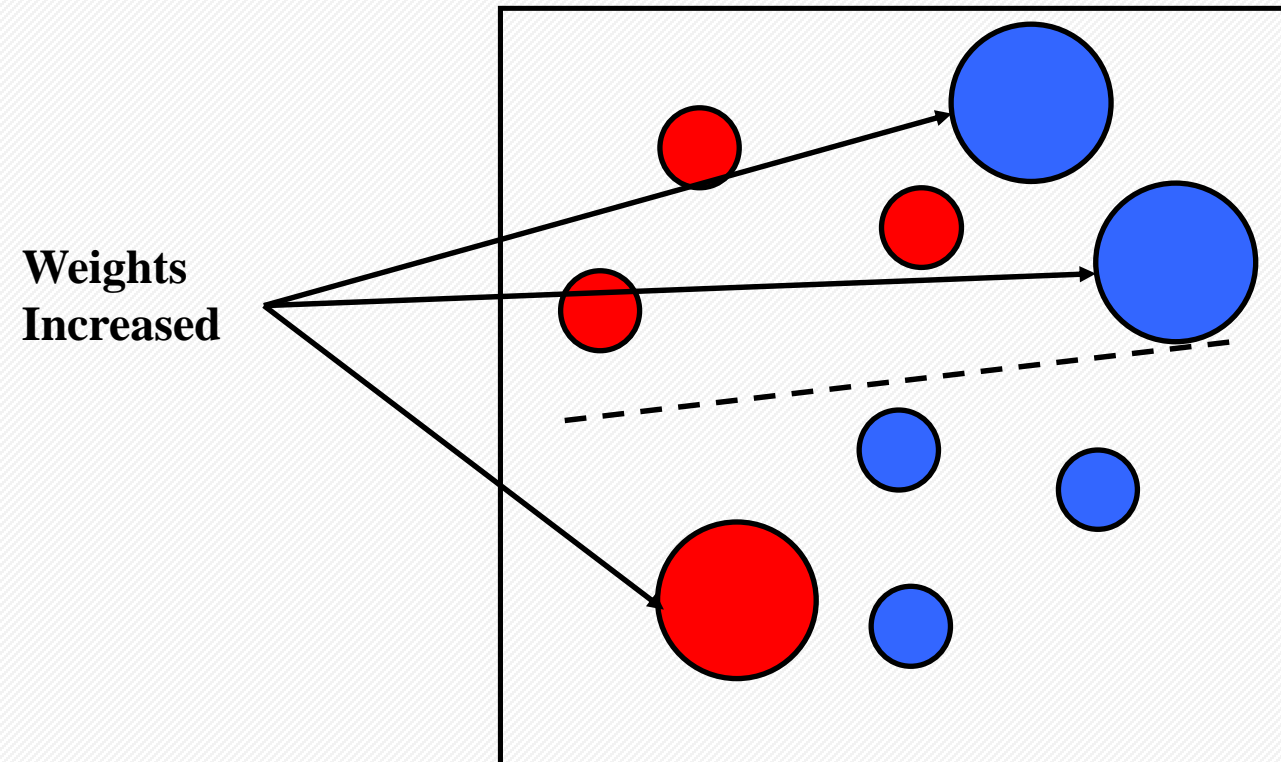
- 当阈值取2.5时，误差率为0.31818。即 $x < 2.5$ 时取 1, $x > 2.5$ 时取 -1, 则数据6、7、8分错，误差率为 $0.10606 * 3 = 0.31818$
- 当阈值取5.5时，误差率最低为0.18184。即 $x < 5.5$ 时取1, $x > 5.5$ 时取 -1, 则数据3、4、5、6、7、8分错，错误率为 $0.16667 * 3 + 0.10606 * 3 = 0.81819 > 0.5$, 故反过来，令 $x > 5.5$ 取 1, $x < 5.5$ 时取 -1, 则数据0、1、2、9分错，误差率为 $0.04546 * 4 = 0.18184$
- 当阈值取8.5时，误差率为0.13638。即 $x < 8.5$ 时取1, $x > 8.5$ 时取 -1, 则数据3、4、5分错，错误率为 $0.04546 * 3 = 0.13638$ (取过，不列入考虑范围)
- 由上面可知，阈值取8.5时，误差率最小，但8.5取过了，所以取5.5。

Boosting illustration

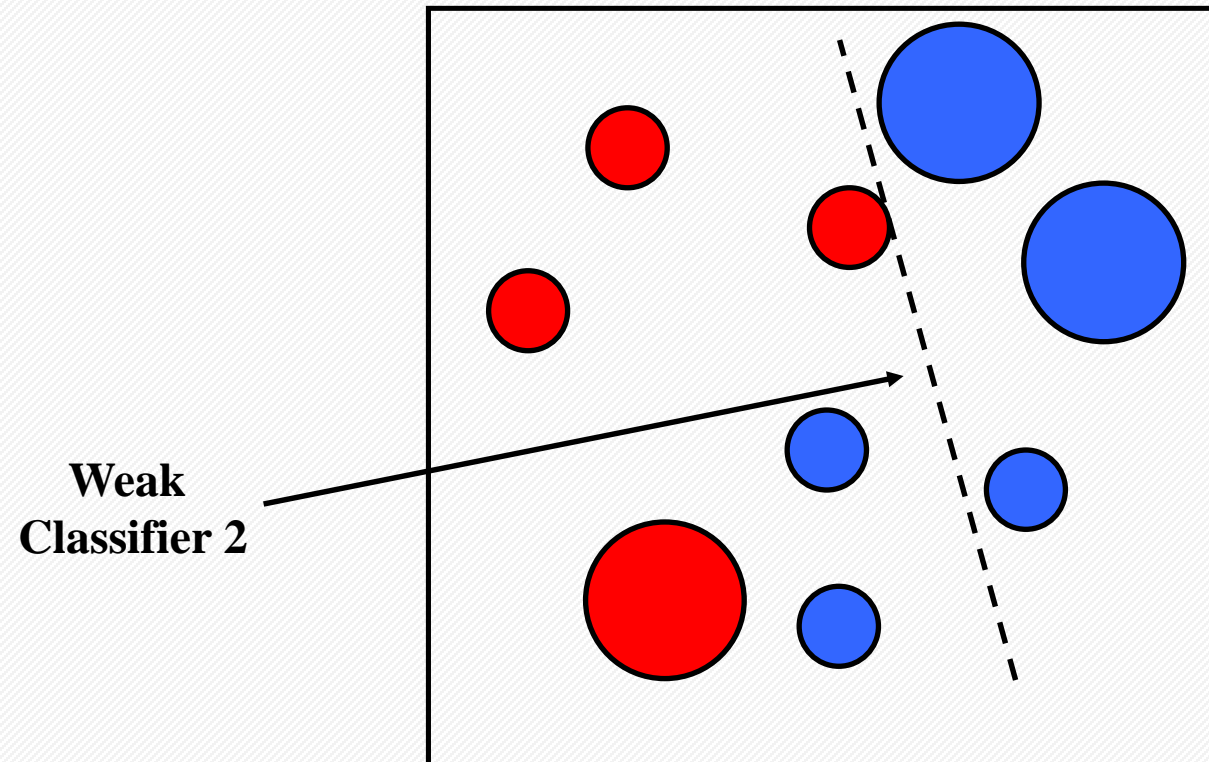
Weak
Classifier 1



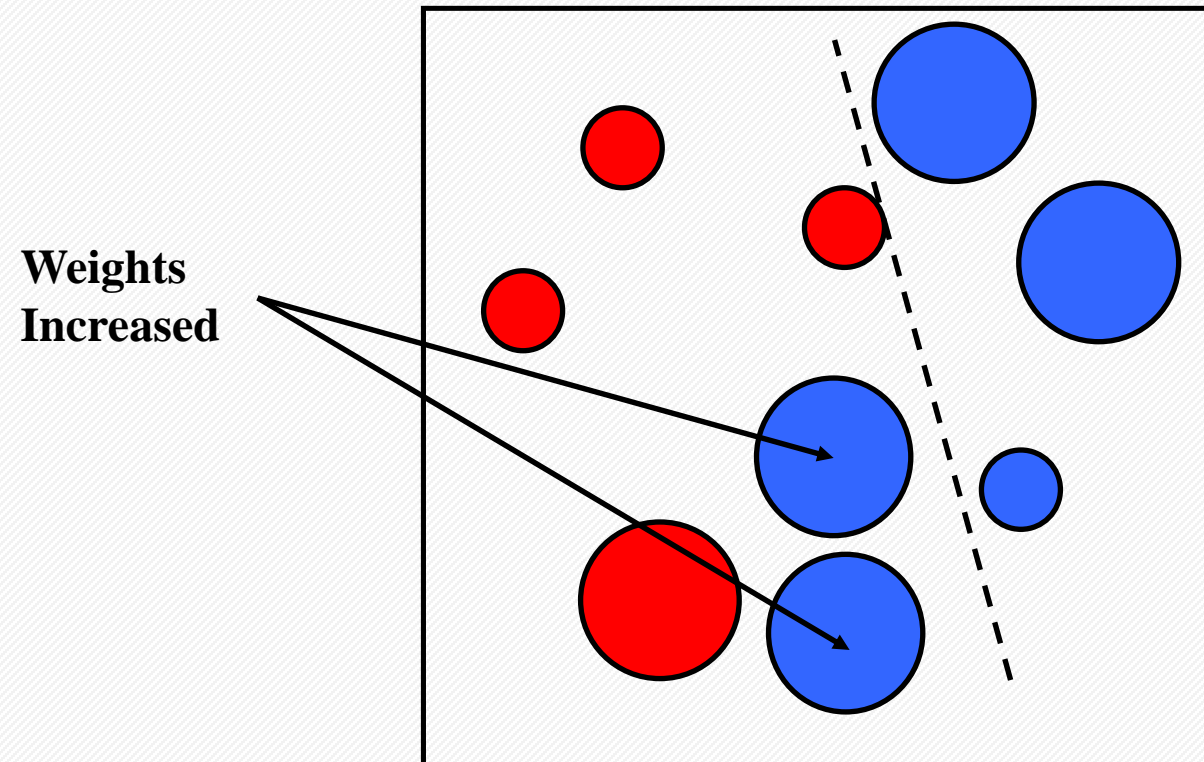
Boosting illustration



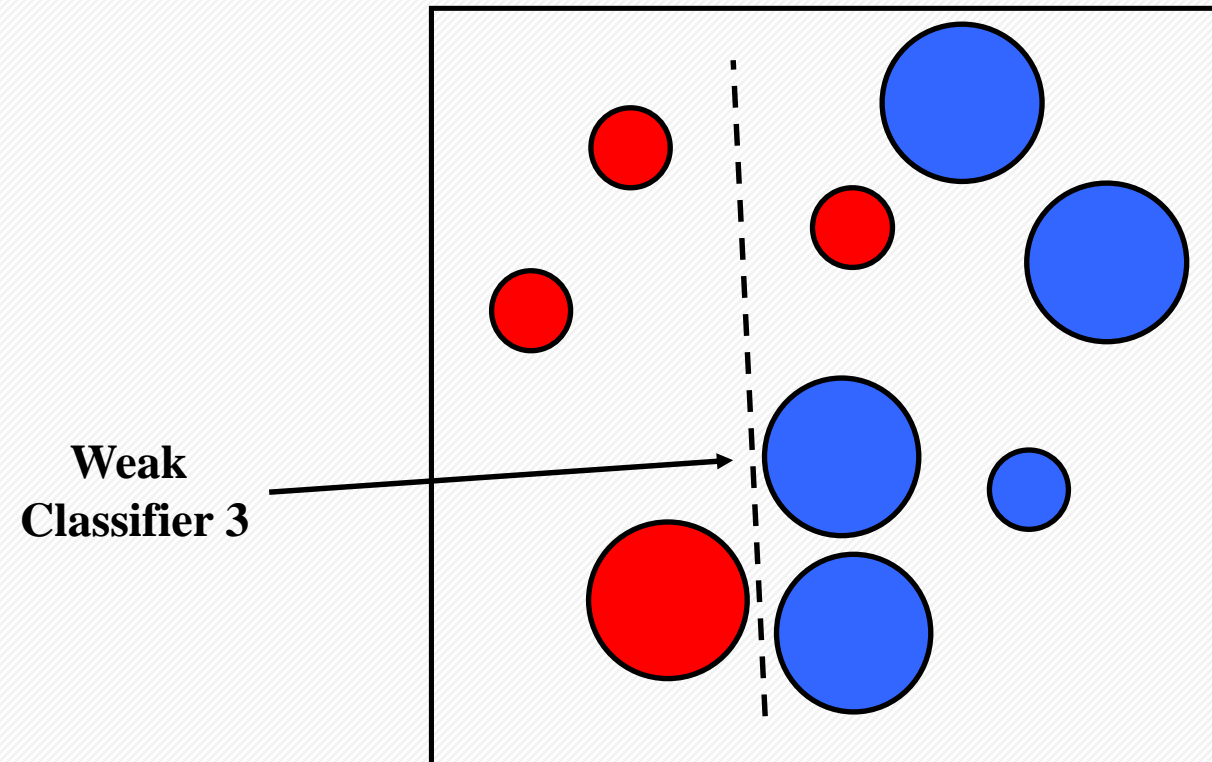
Boosting illustration



Boosting illustration

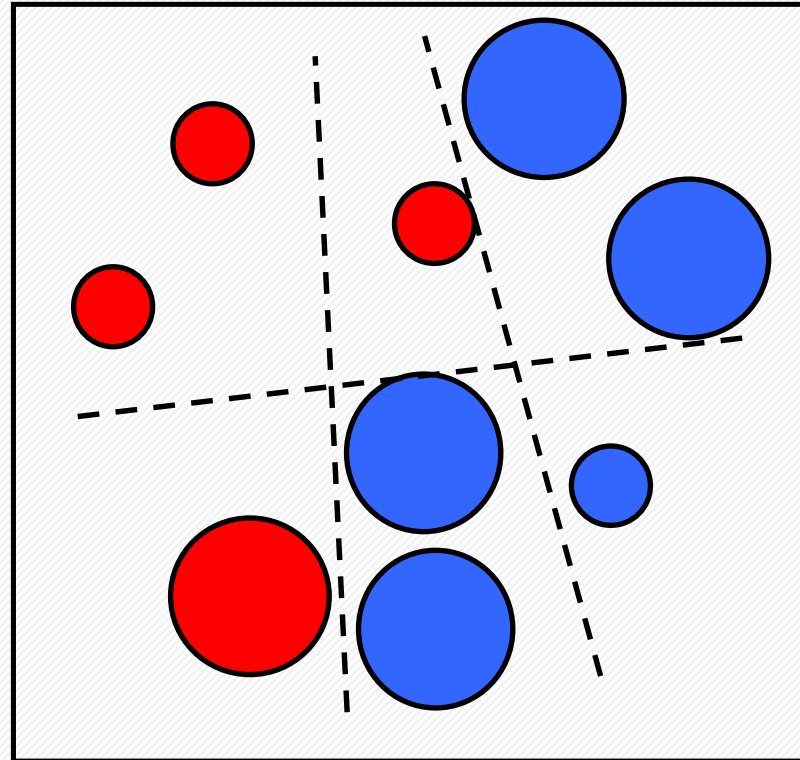


Boosting illustration



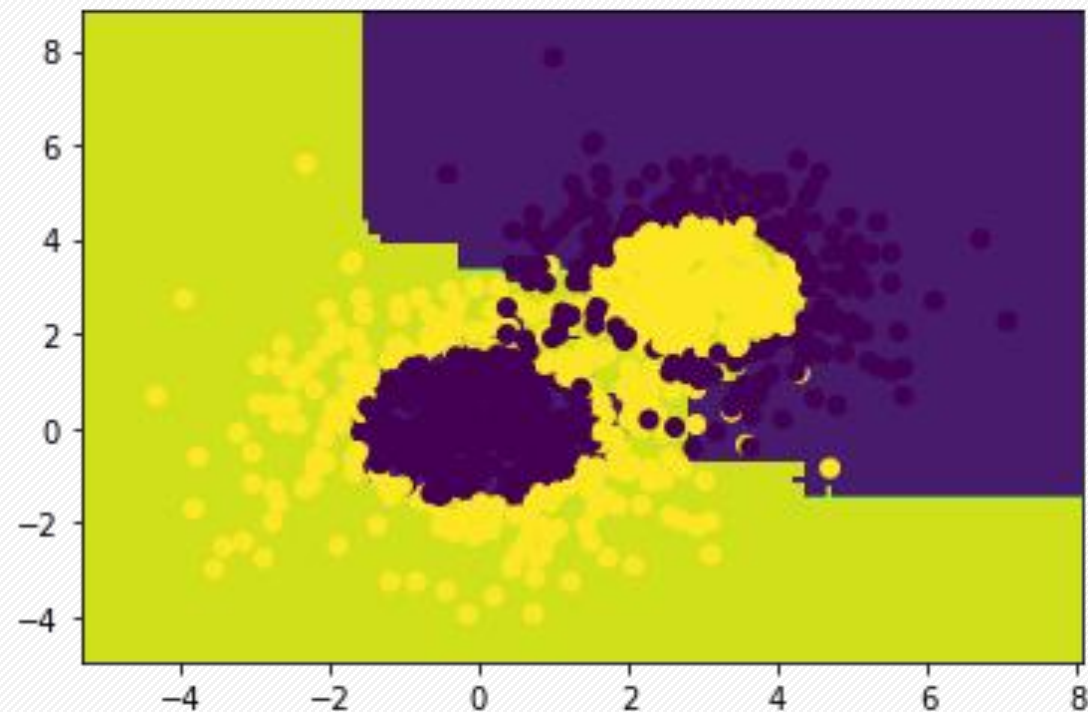
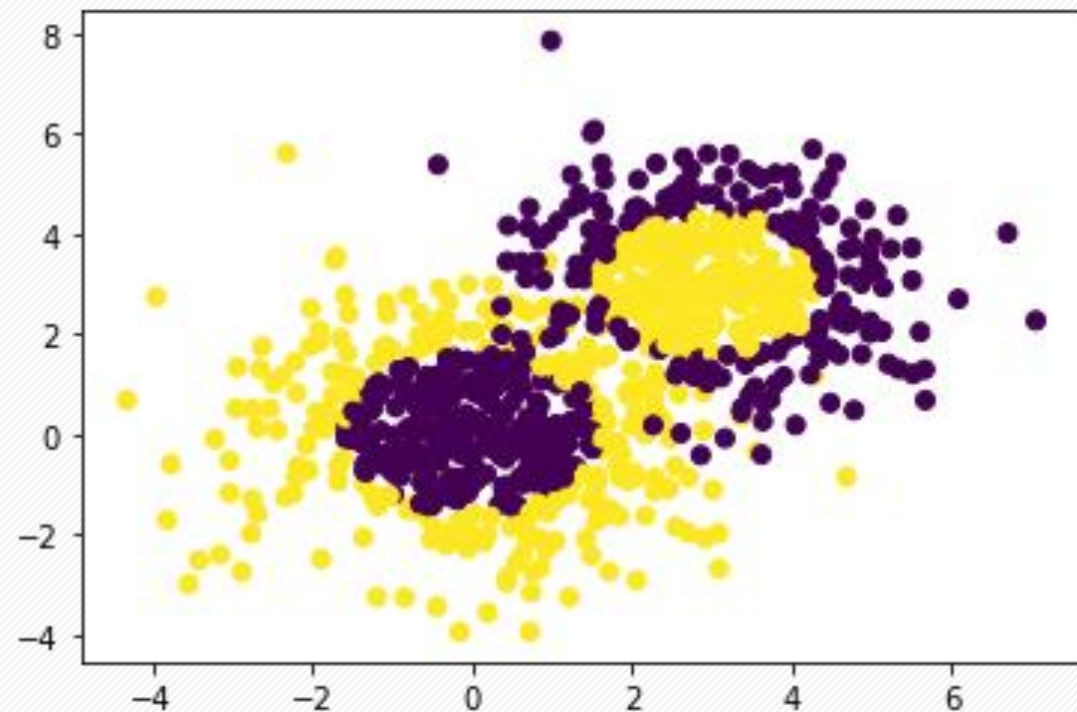
Boosting illustration

**Final classifier is
a combination of weak
classifiers**





例8.2



8.3 Bagging与随机森林

8.3.1 Bagging

- 给定含 m 个样本的数据集，先随机取出一个样本放入采样集中，再把该样本放回初始数据集。使得下次采样时该样本仍有可能被选中，经过 m 次随机采样操作，得到包含 m 个样本的采样集。
- 新数据集和原数据集的大小相等。
- 采样 T 个含 m 个训练样本的采样集，然后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。
- Bagging通常对分类任务使用简单投票法，对回归任务使用简单平均法。

输入： 训练数据集 $D = \{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_N, y_N)\}$,
基学习算法 \mathcal{L} ; 训练轮数 T

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

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

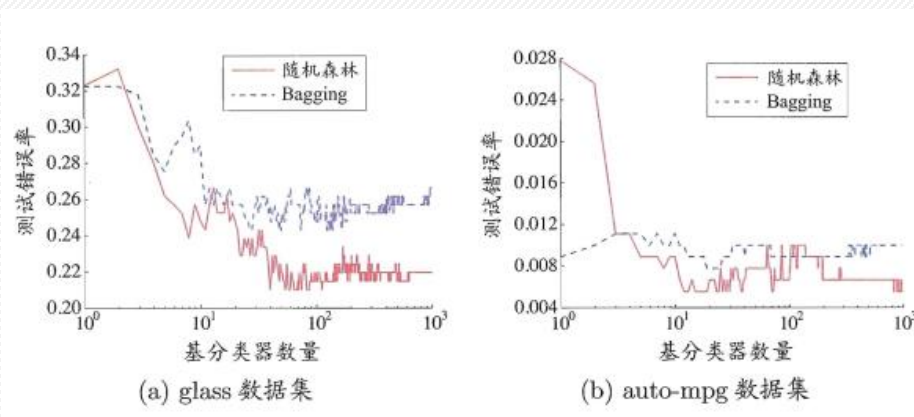


8.3.2 随机森林(Random Forest,RF)

Bagging的一个扩展变体。RF在以**决策树**为基础学习器构建Bagging集成的基础上，进一步在决策树的训练过程中引入**随机属性选择**。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。一般情况下，参数控制了随机性的引入程度。推荐值 $k = \log_2 d$ 。

优点：随机森林简单、容易实现、计算开销小。

被誉为“**代表集成学习技术水平的方法**”。



8.4 结合策略

学习器结合可能会从三个方面带来好处[Dietterich, 2000]:

(1) **统计**: 由于学习任务的假设空间往往很大, 可能有多个假设在训练集上达到同等性能, 此时若使用单学习器可能因误选而导致泛化性能不佳, **结合多个学习器则会减小这一风险**;

(2) **计算**: 学习算法往往会陷入局部极小, 有的局部极小点所对应的泛化性能可能很糟糕, 而通过多次运行之后进行结合, 可**降低陷入糟糕局部极小点的风险**;

(3) **表示**: 某些学习任务的真实假设可能不在当前学习算法所考虑的假设空间中, 此时若使用单学习器则肯定无效, 而通过结合多个学习器, 由于相应的**假设空间有所扩大**, 有可能学得更好的近似。



8.4 结合策略

常用结合方法

- **平均法**
 - 简单平均法
 - 加权平均法
- **投票法**
 - 绝对多数投票法
 - 相对多数投票法
 - 加权投票法
- **学习法**



8.4.1 平均法

•简单平均法 (simple averaging)

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

•加权平均法 (weighted averaging)

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

$$\omega_i \geq 0, \text{ 且 } \sum_{i=1}^T \omega_i = 1$$

其中，加权平均法的**权重**一般从训练数据中学习得到。但由于噪声干扰和过拟合问题，加权平均法也未必优于简单平均法[Xu et al, 1992; Ho et al, 1994; Kittler et al., 1998]。

推荐：个体学习器性能相差较大时，宜用加权平均法，否则宜用简单平均法。

8.4.2 投票法

将 h_i 在样本 x 上的预测输出表示为一个 N 维向量 $(h_i^1(x), h_i^2(x), \dots, h_i^N(x))$, 其中 $h_i^j(x)$ 是 h_i 在类别标记 c_j 上的输出。

• 绝对多数投票法 (majority voting)

$$H(x) = \begin{cases} c_j, & \text{if } \sum_{i=1}^T h_i^j(x) > 0.5 \sum_{k=1}^N \sum_{i=1}^T h_i^k(x) \\ \text{reject}, & \text{otherwise} \end{cases}$$

• 相对多数投票法 (plurality voting)

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

• 加权投票法 (weighted voting)

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



8.4.3 学习法

Stacking 算法

输入: 训练集 $D = \{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_m, y_m)\}$;
初级学习算法 $\mathcal{L}_1, \mathcal{L}_2, \dots, \mathcal{L}_T$;
次级学习算法 \mathcal{L} .

过程:

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

使用初级学习算法产生初级学习器

```
4:  $D' = \emptyset$ ;
```

```
5: for  $i = 1, 2, \dots, m$  do  
6:   for  $t = 1, 2, \dots, T$  do  
7:      $z_{it} = h_t(\mathbf{x}_i)$ ;  
8:   end for  
9:    $D' = D' \cup ((z_{i1}, z_{i2}, \dots, z_{iT}), y_i)$ ;  
10: end for
```

生成次级训练集

```
11:  $h' = \mathcal{L}(D')$ ;
```

输出: $H(\mathbf{x}) = h'(h_1(\mathbf{x}), h_2(\mathbf{x}), \dots, h_T(\mathbf{x}))$



8.5 多样性

8.5.1 误差-分歧分解

假定用个体学习器 h_1, h_2, \dots, h_T 通过加权平均法结合产生集成来完成回归学习任务 $f: R^d \mapsto R$ 。

$$H(\mathbf{x}) = \sum_{i=1}^T \omega_i h_i(\mathbf{x}), \text{ 其中 } \omega_i \geq 0 \text{ 是学习器权重, } \sum_{i=1}^T \omega_i = 1$$

1、分歧

个体学习器 h_i 的分歧

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

集成的分歧

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



8.5.1 误差-分歧分解

2、个体学习器 h_i 和集成 H 的平方误差

个体学习器 h_i 的平方误差：

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

集成 H 的平方误差：

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

个体学习器误差的加权均值：

$$\bar{E}(h|\mathbf{x}) = \sum_{i=1}^T \omega_i E(h_i|\mathbf{x})$$



8.5.1 误差-分歧分解

4、误差-分歧分解 (error-ambiguity decomposition)

step 1: 将平方误差带入集成的分歧中, 得:

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

step 2: 令 $p(x)$ 表示样本概率密度, 在全样本上有 (对step 1式两边对 x 积分) :

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

step 3: 个体学习器 h_i 在全样本上的泛化误差:

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



8.5.1 误差-分歧分解

个体学习器 h_i 泛化误差的加权均值:

$$\bar{E} = \sum_{i=1}^T \omega_i E_i$$

个体学习器 h_i 在全样本上的分歧项:

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

个体学习器 h_i 的加权分歧值:

$$\bar{A} = \sum_{i=1}^T \omega_i A_i$$

集成的泛化误差:

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

集成的泛化误差 E (误差-分歧分解):

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

上面这个分析首先由[Krogh and Vedelsby, 1995]给出, 称为误差-分歧分解。

现实任务中很难直接对该式进行优化



8.5.2 多样性度量

给定数据集 $D = \{(x_1, y_1), (x_2, y_2), \dots, (x_m, y_m)\}$, 对于二分类任务, $y_i \in \{-1, 1\}$, 分类器 h_i 与 h_j 的预测结果列联表 (contingency table) 为:

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

• 不合度量 (disagreement measure)

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

• 相关系数 (correlation coefficient)

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

8.5.2 多样性度量

•Q-统计量 (Q-statistic)

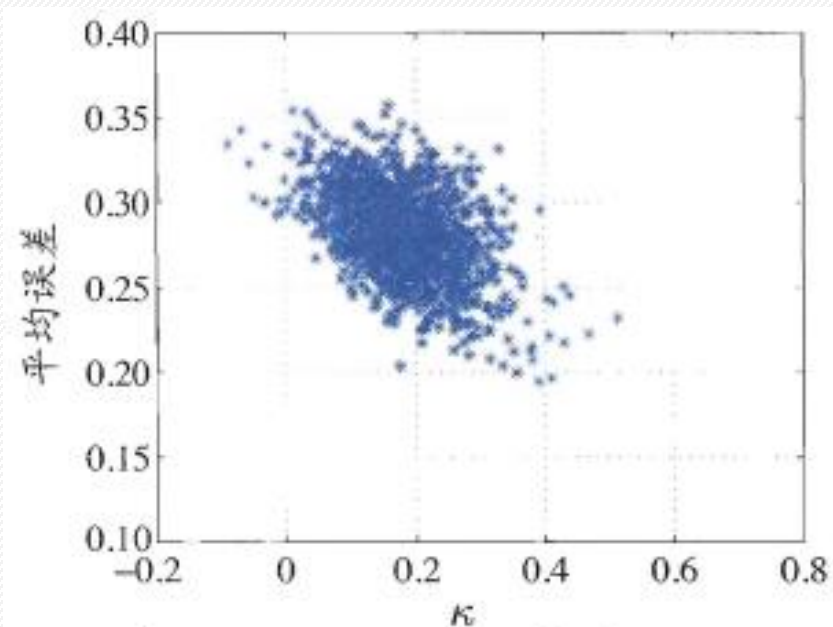
$$Q_{ij} = \frac{ad - bc}{ad + bc}$$

•k-统计量 (k-statistic)

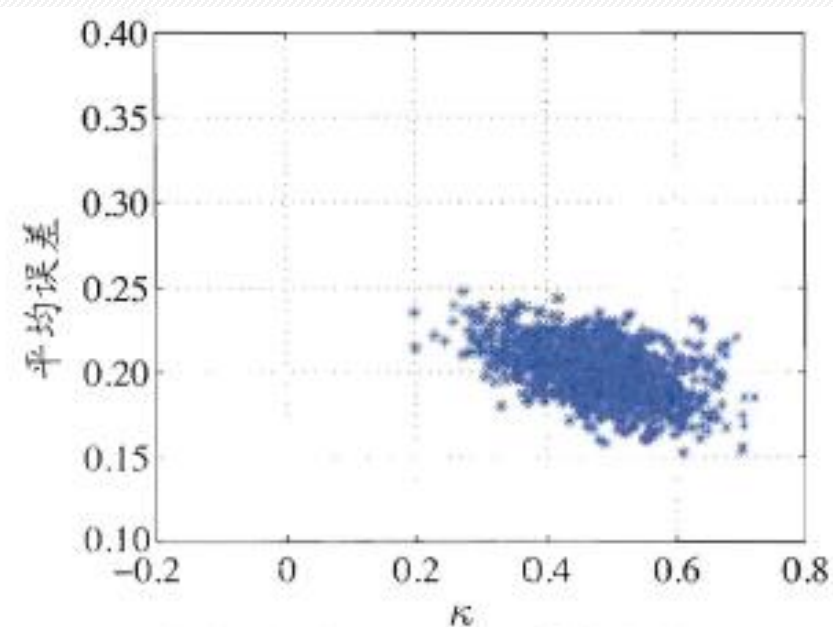
$$k = \frac{p_1 - p_2}{1 - p_2}$$

$$p_1 = \frac{a + d}{m}$$

$$p_2 = \frac{(a + b)(a + c) + (c + d)(b + d)}{m^2}$$



(a) AdaBoost 集成



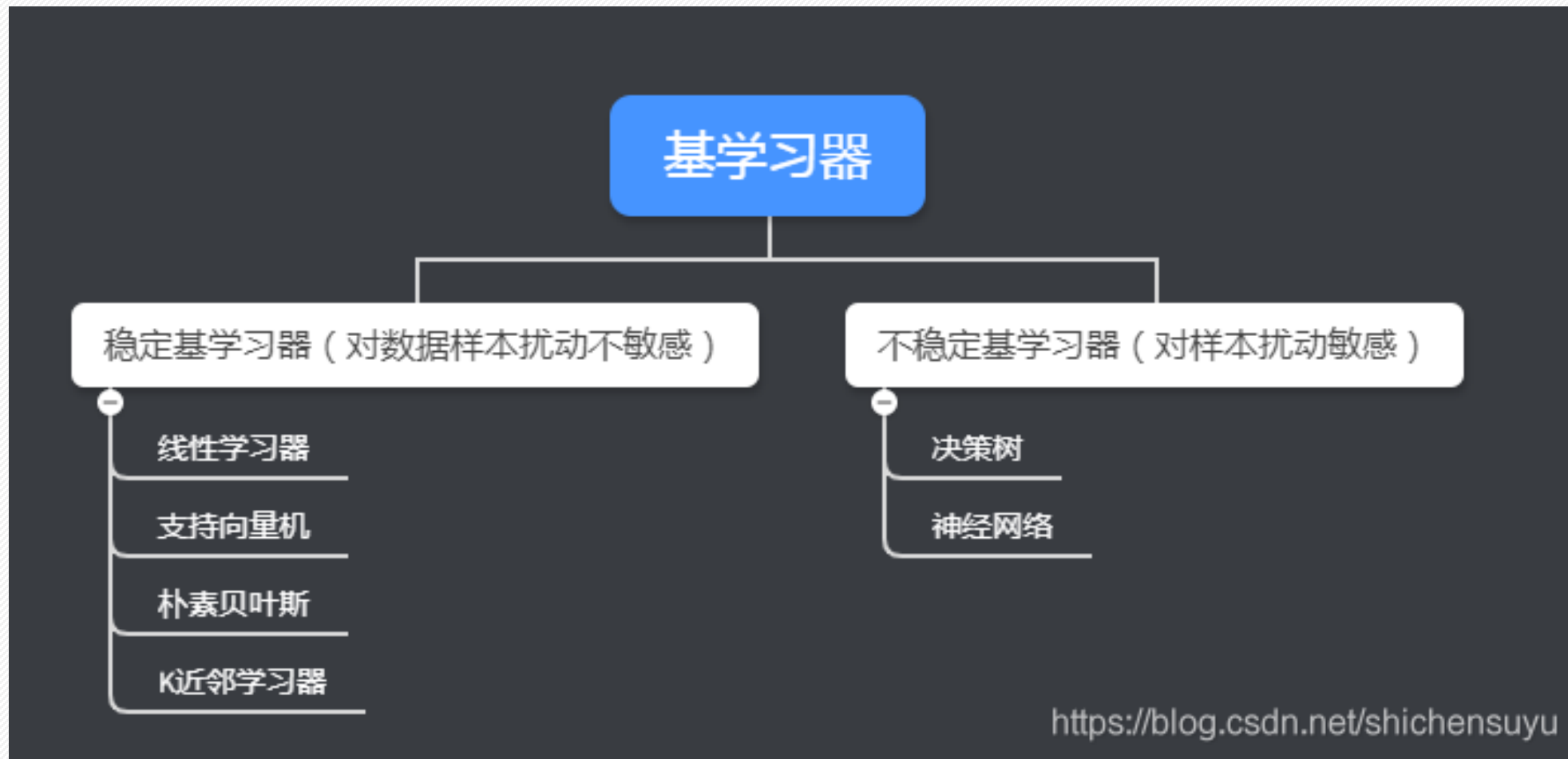
(b) Bagging 集成



8.5.3 多样性增强

- **数据样本扰动（适合不稳定基学习器，如决策树、神经网络等）**

通常基于采样法，如在Bagging中使用自助采样，在AdaBoost采用序列采样





8.5.3 多样性增强

- **输入属性扰动 (适合包含大量冗余属性的数据)**

不同的子空间 (即属性子集) 提供了观察数据的不同视角。显然, 从不同的属性子集训练出的个体学习器必然有所不同。

随机子空间(random subspace)算法采用输入属性扰动: **从初始属性集中抽取**
出若干个属性子集, 再基于每个属性子集训练一个基学习器。

d' 小于初始属性数 d .

\mathcal{F}_t 包含 d' 个随机选取
的属性, D_t 仅保留 \mathcal{F}_t 中
的属性.

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} ;
基学习器数 T ;
子空间属性数 d' .

过程:

```
1: for  $t = 1, 2, \dots, T$  do
2:    $\mathcal{F}_t = \text{RS}(D, d')$ 
3:    $D_t = \text{Map}_{\mathcal{F}_t}(D)$ 
4:    $h_t = \mathcal{L}(D_t)$ 
5: end for
```

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

图 8.11 随机子空间算法 cdn.net/shichensuyu



8.5.3 多样性增强

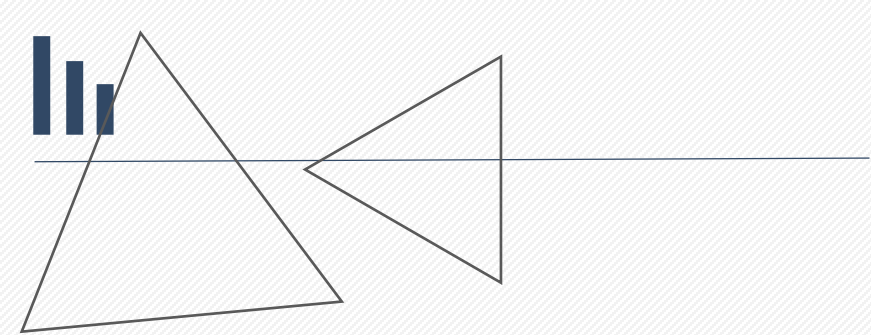
- **输出表示扰动**

对输出表示进行操纵可以增强多样性对训练样本的类标记稍作变动，如翻转法：随机改变一些训练样本的标记对输出表示进行转化，如输出调制法：将分类输出转化为回归输出后构建个体学习器将原任务拆解为多个可同时求解的子任务，如ECOC法：利用纠错输出码将多分类任务拆解为一系列二分类来训练基学习器

- **算法参数扰动**

通过随机设置不同的参数，往往可产生差别较大的个体学习器。

例如负相关法：显式地通过正则化项来强制个体神经网络使用不同的参数。使用单一学习器时通常使用交叉验证等方法确定参数值，这事实上已经使用了不同参数训练出多个学习器，只不过最终选取其中一个学习器进行使用，而集成学习则相当于把这些学习器都利用起来。



The end

