



大数据导论

Introduction to Big Data



第10讲: 集成学习

叶允明

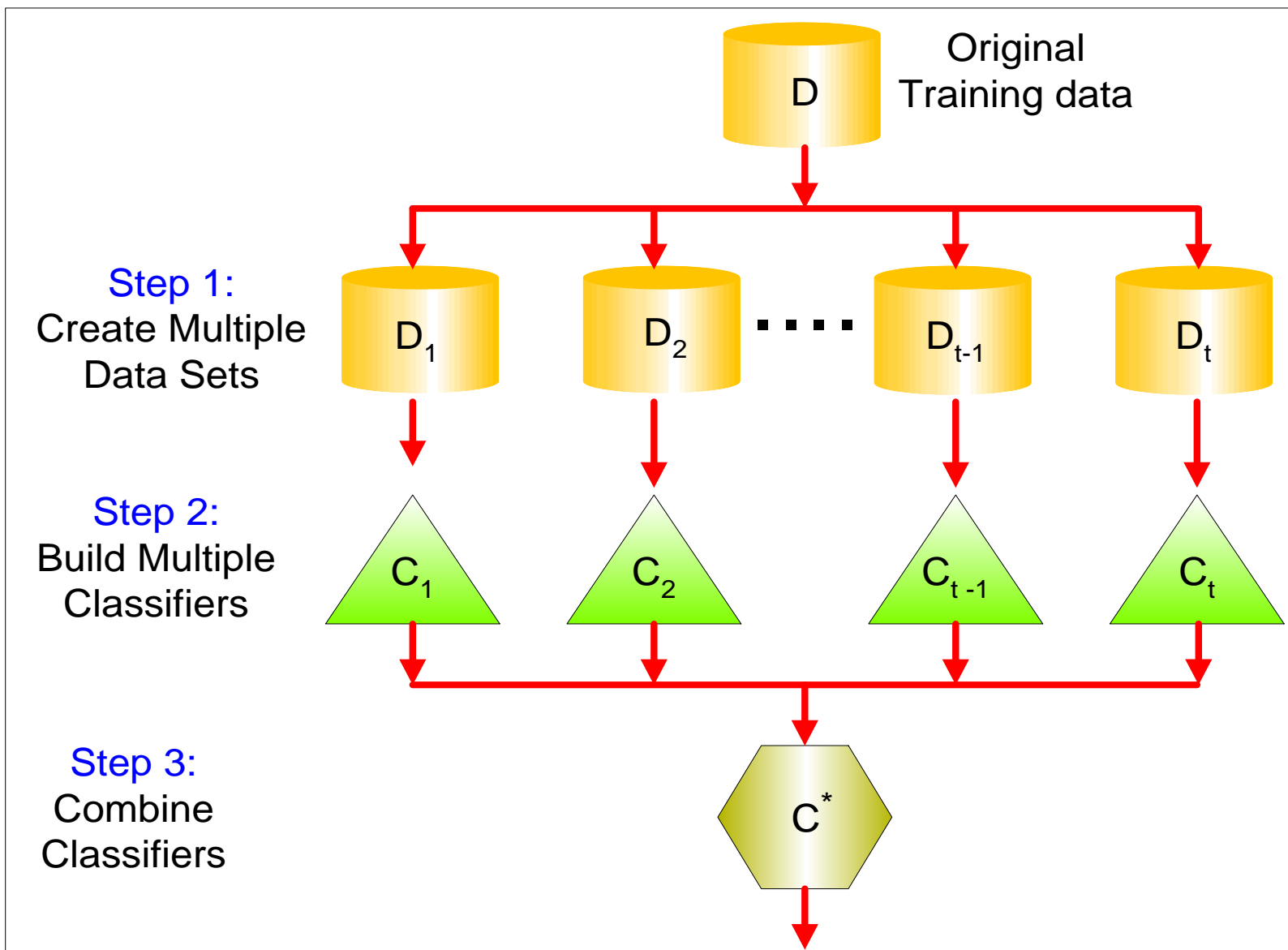
计算机科学与技术学院

哈尔滨工业大学 (深圳)

基于集成学习的分类方法

- 集成学习 (ensemble learning)
- 从训练数据构建分类器集合
- 综合多个分类器的预测结果来预测新实例的分类预测结果

总体思路



为什么会起作用？

- 假设有25个基础分类器
 - 每个分类器的错误率是 $\varepsilon = 0.35$
 - 假设每个分类器之间是相互独立的
 - 那么这个集成分类器预测错误的概率是：

$$\sum_{i=13}^{25} \binom{25}{i} \varepsilon^i (1 - \varepsilon)^{25-i} = 0.06$$

集成学习的关键问题

- 强度 (strength)
- 相关性 (correlation)
- Tradeoff between strength and correlation

集成方法的例子

- 如何生成一个集成分类器？
 - 装袋 (Bagging)
 - 提升 (Boosting)
 - 梯度提升 (Gradient Boosting)
 - 堆叠 (Stacking)

装袋 (Bagging)

- 有放回抽样

Original Data	1	2	3	4	5	6	7	8	9	10
Bagging (Round 1)	7	8	10	8	2	5	10	10	5	9
Bagging (Round 2)	1	4	9	1	2	3	2	7	3	2
Bagging (Round 3)	1	8	5	10	5	5	9	6	3	7

- 为每个自助样本建立分类器
- 每个样本被抽取的概率是相等的

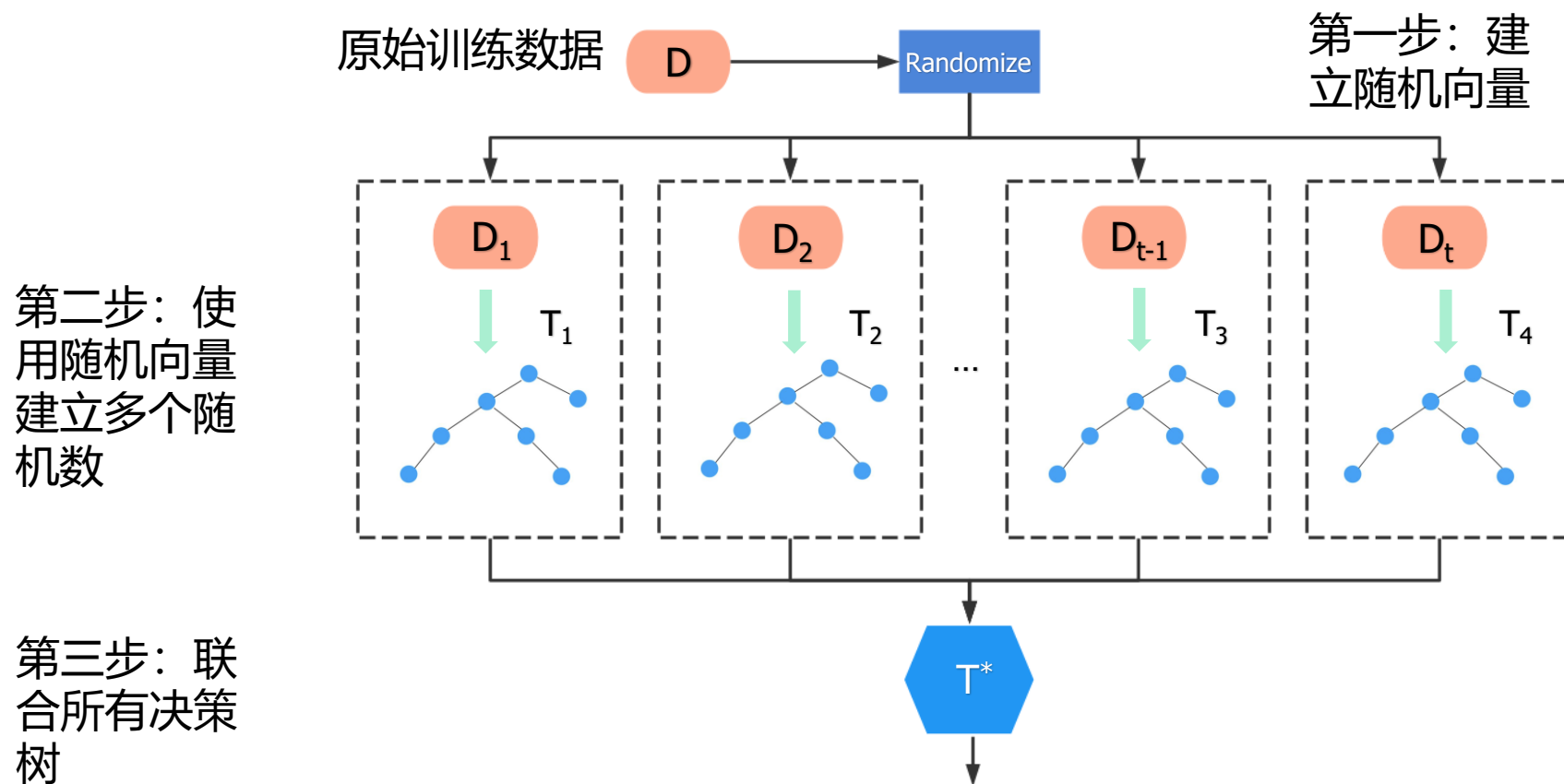
随机森林

- 专门为决策树分类器设计的集成方法
- 随机森林长出许多分类树 (名字由来)
- 未修剪的决策树集合
- 每个基础分类器分类一个 “新的” 向量
- 森林选择表决投票最多的分类结果

随机森林

- 随机性的两个来源介绍: “装袋” 和 “随机输入向量”
 - 每棵树都是使用训练数据的bootstrapping样本来生长的
 - 在每个节点上, 最佳的分裂是从 m_{try} 个变量中选择, 而不是所有变量

随机森林



提升 (Boosting)

- 更注重之前错误分类的实例以自适应地改变训练数据分布的迭代过程
 - 最初，所有N个实例都会分配相同的权重
 - 跟装袋不同的是，这些权重会在每一轮提升结束之后发生变化
 - 最终的分类器是基于弱分类器的权重表决的

Original Data	1	2	3	4	5	6	7	8	9	10
Boosting (Round 1)	7	3	2	8	7	9	4	10	6	3
Boosting (Round 2)	5	4	9	4	2	5	1	7	4	2
Boosting (Round 3)	4	4	8	10	4	5	4	6	3	4

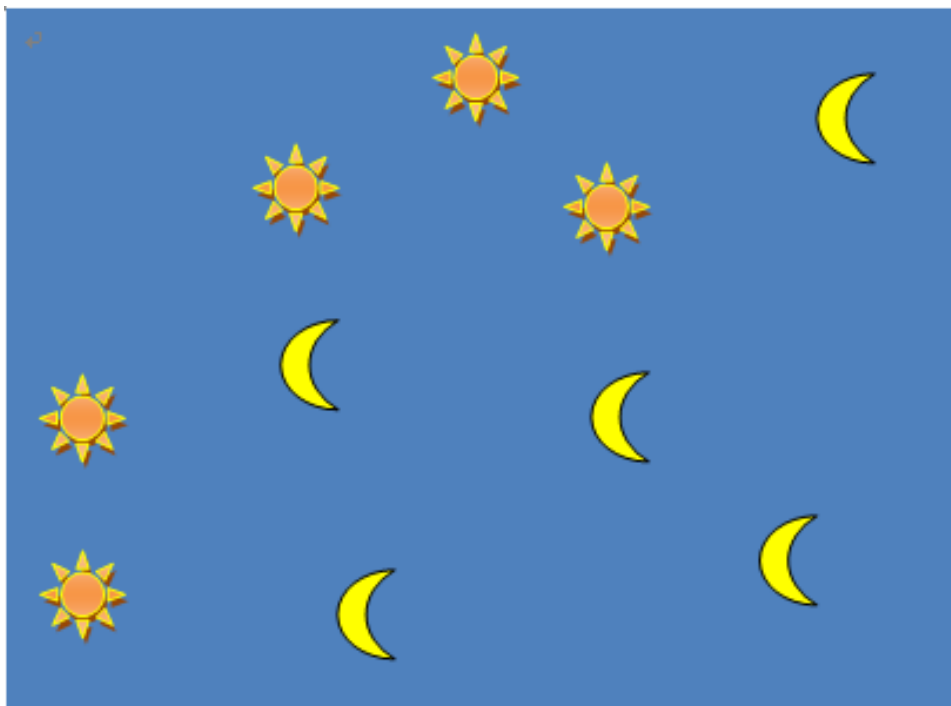
Adaboost (Freund and Schapire, 1997)

- 给定数据集D, 包含d个类标记的元组 $(\mathbf{X}_1, y_1), \dots, (\mathbf{X}_d, y_d)$
- 初始情况下, 对每个训练元组赋予相等的权重 $1/d$
- 执行k轮算法产生k个基分类器。在第 i 轮:
 - 从D中元组进行有放回抽样, 形成相同大小的训练集 D_i
 - 每个元组被选中的机会由它的权重决定
 - 从训练集 D_i 中导出分类器 M_i
 - 使用 D_i 作为验证集计算 M_i 的错误率
 - 如果元组被不正确地分类, 那么它的权重增加; 否则权重减少
- 错误率: $err(\mathbf{X}_j)$ 是元组 \mathbf{X}_j 的误分类误差; 分类器 M_i 的错误率是每个元组错误率的加权和:

$$error(M_i) = \sum_j^d w_j \times err(\mathbf{X}_j)$$

- 分类器 M_i 的表权重为: $\log \frac{1 - error(M_i)}{error(M_i)}$

Adaboost:例子



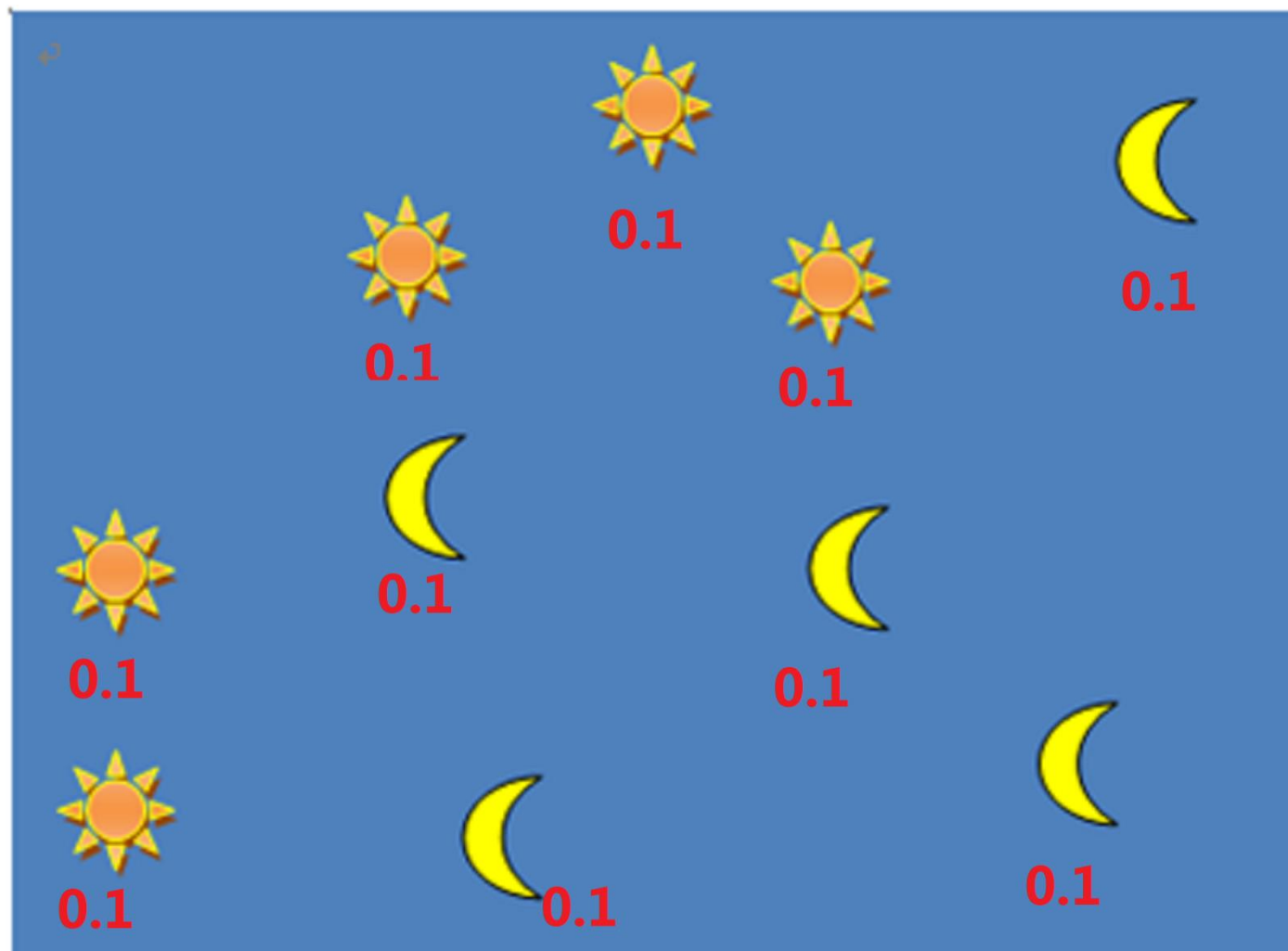
☀ 和 🌙 分别代表两个类别

在此过程中，我们使用水平或垂直线作为分类器进行分类。

Adaboost:例子

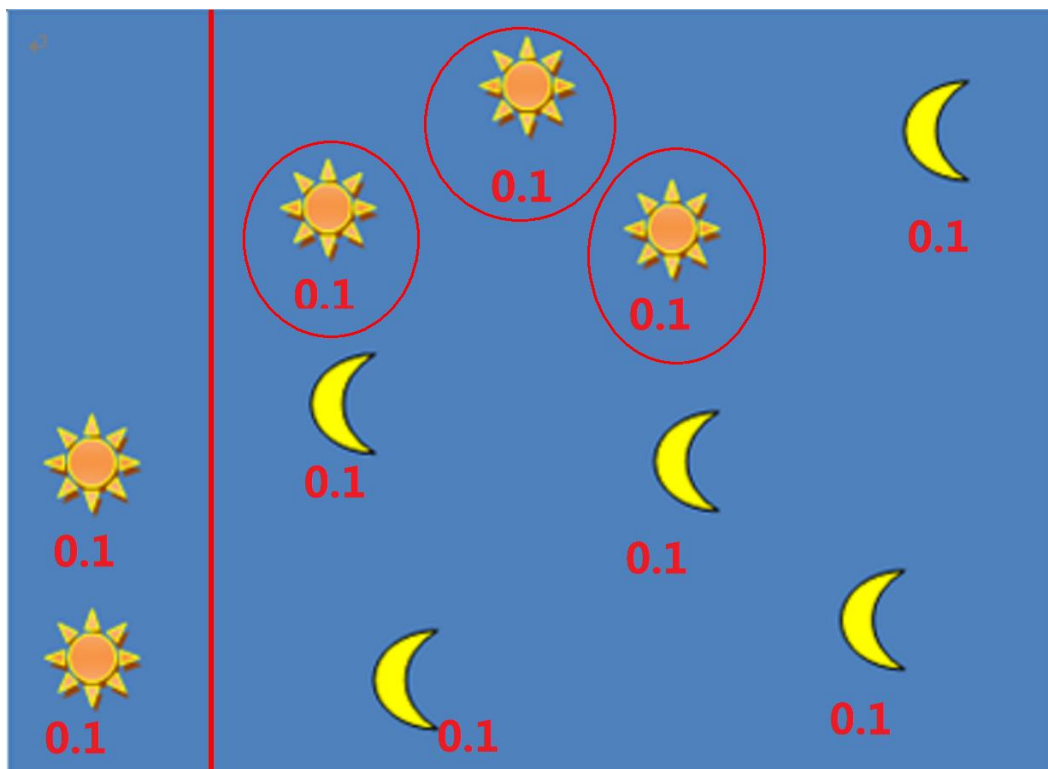
初始化元组的权重:

$$w_{\text{init}} = \frac{1}{D} = \frac{1}{10} = 0.1$$



Adaboost:例子

第一轮: M1



M1的错误率:

$$\begin{aligned} error(M_1) &= \sum_{j=1}^{10} w_j \times err(X_j) \\ &= 0.1 + 0.1 + 0.1 \\ &= 0.3 \end{aligned}$$

M1的表决权重:

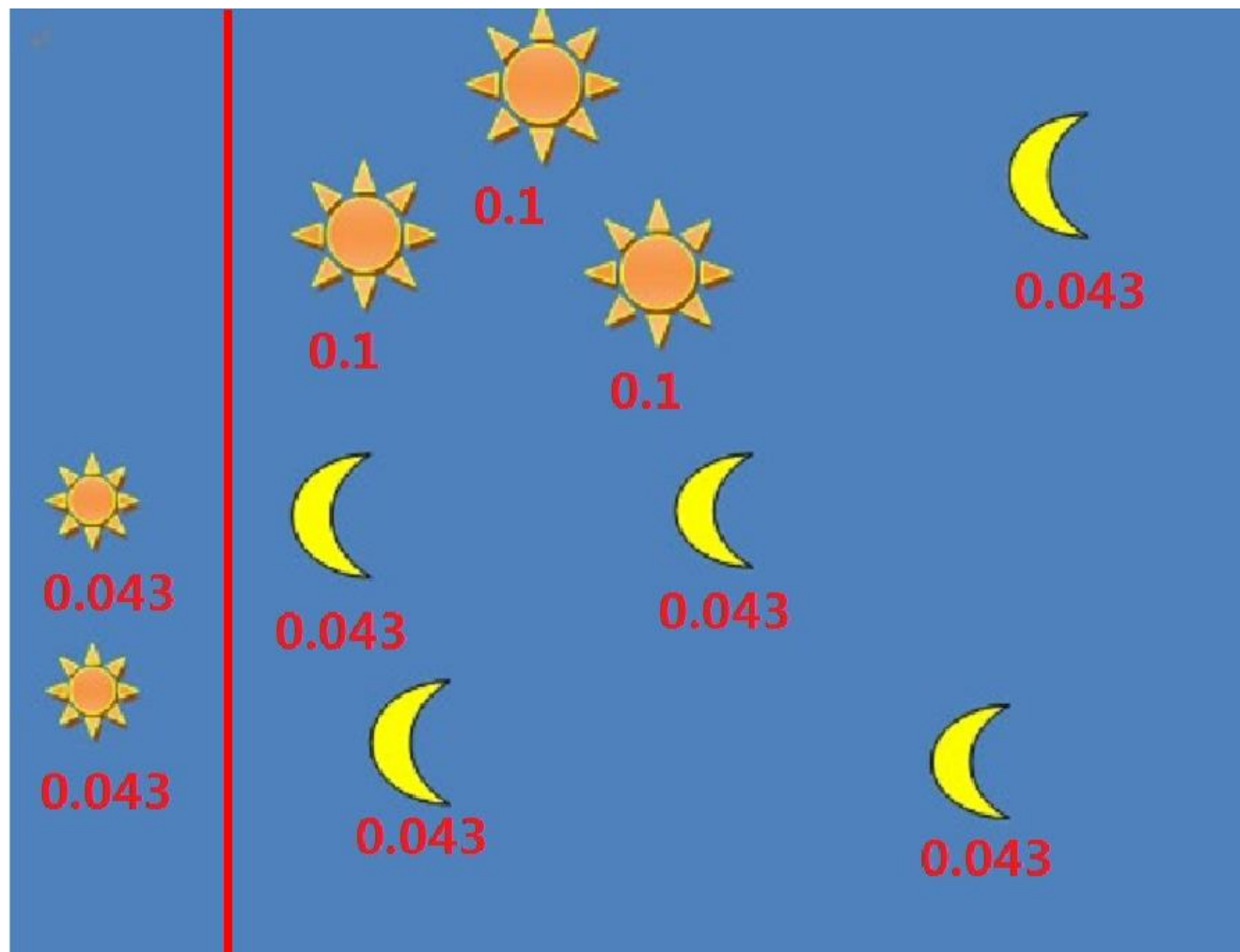
$$\begin{aligned} W_1 &= \log \frac{1 - error(M_1)}{error(M_1)} \\ &= 0.477 \end{aligned}$$

Adaboost:例子

第一轮: M1

更新被正确分类的元组的权重:

$$0.1 * \frac{error(M_1)}{1 - error(M_1)} \\ = 0.043$$



Adaboost:例子

归一化:

旧的权重之和:

$$0.1 * 10 = 1.0$$

新的权重之和:

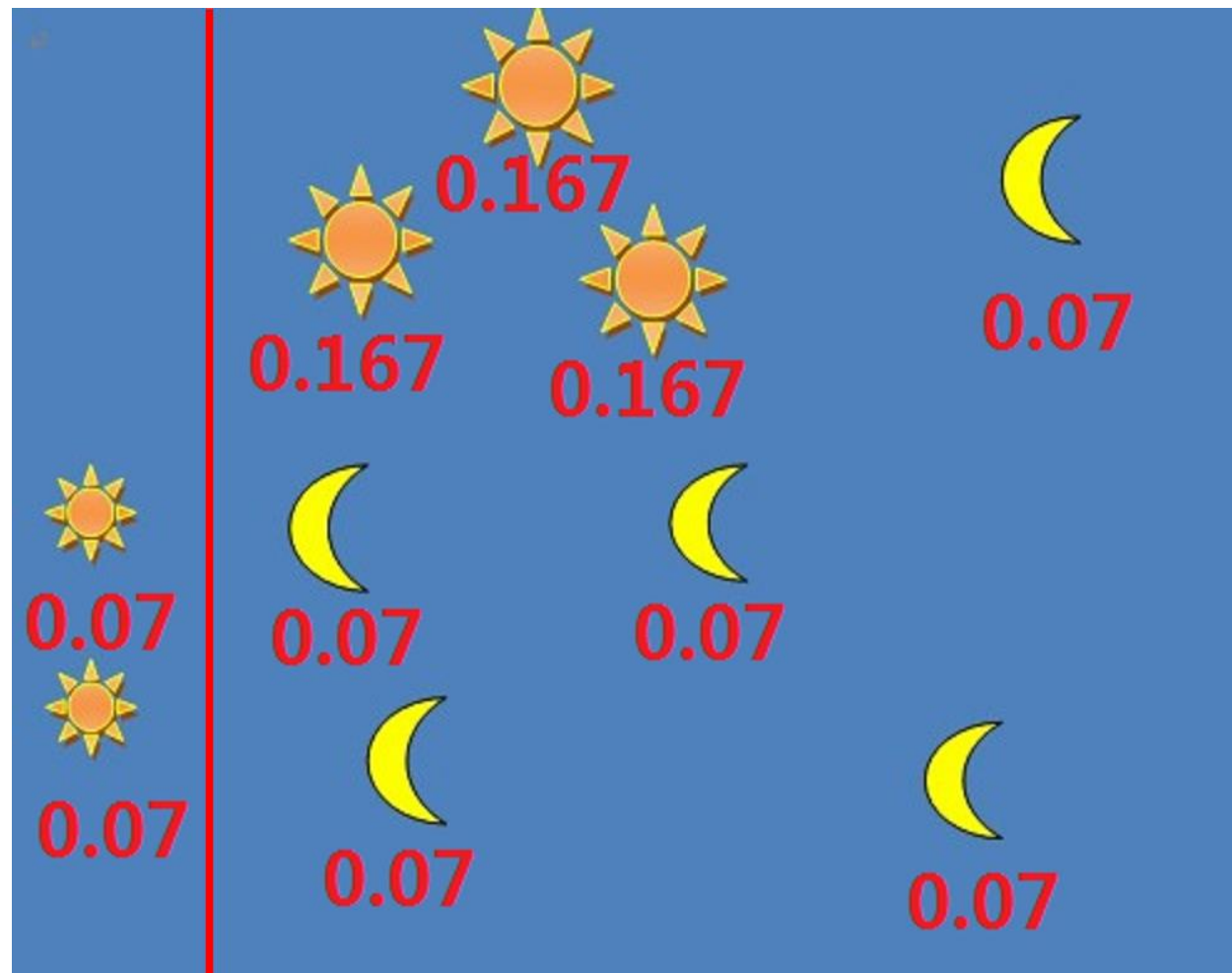
$$0.043 * 7 + 0.1 * 3 = 0.6$$

错误地分类:

$$0.1 * \frac{1}{0.6} = 0.167$$

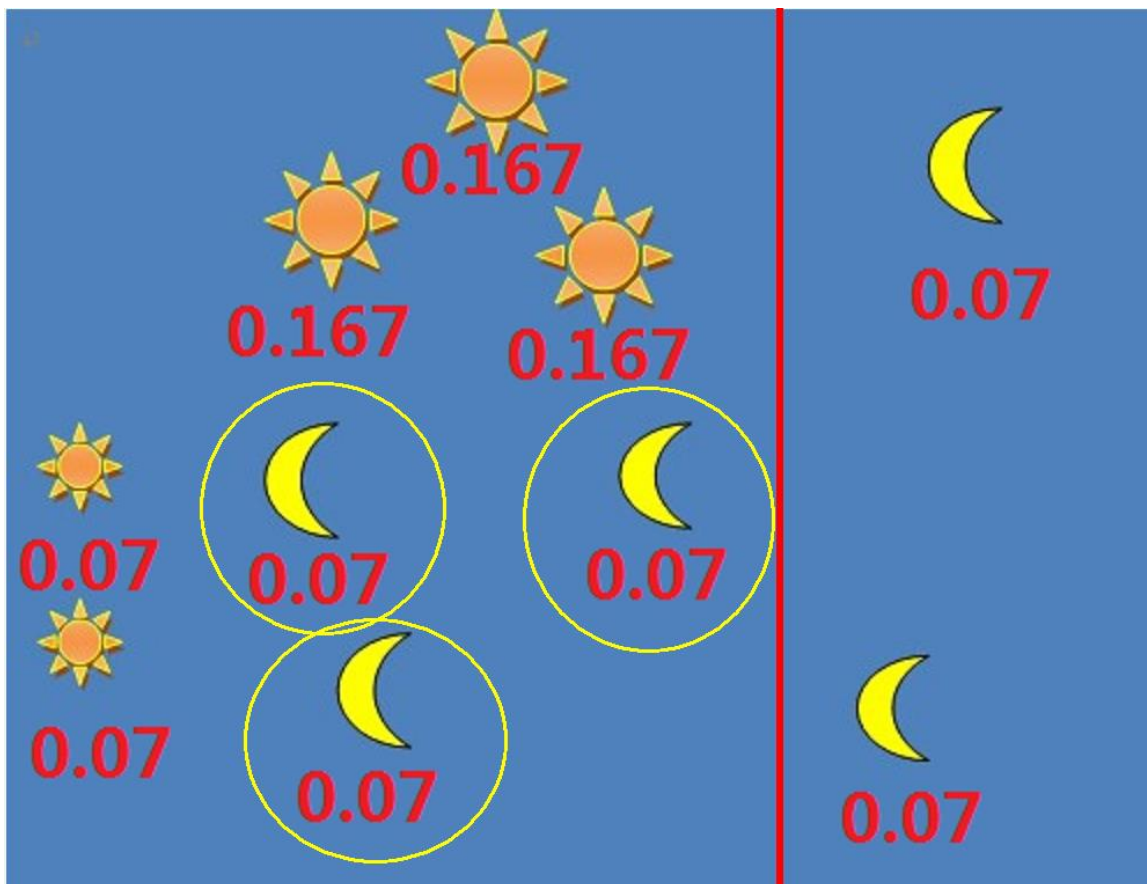
正确地分类:

$$0.043 * \frac{1}{0.6} = 0.07$$



Adaboost:例子

第二轮: M2



M2的错误率:

$$\begin{aligned} error(M_2) &= \sum_{j=1}^{10} w_j \times err(X_j) \\ &= 0.07 + 0.07 + 0.07 \\ &= 0.21 \end{aligned}$$

M2的表决权重:

$$\begin{aligned} W_2 &= \log \frac{1 - error(M_2)}{error(M_2)} \\ &= 0.575 \end{aligned}$$

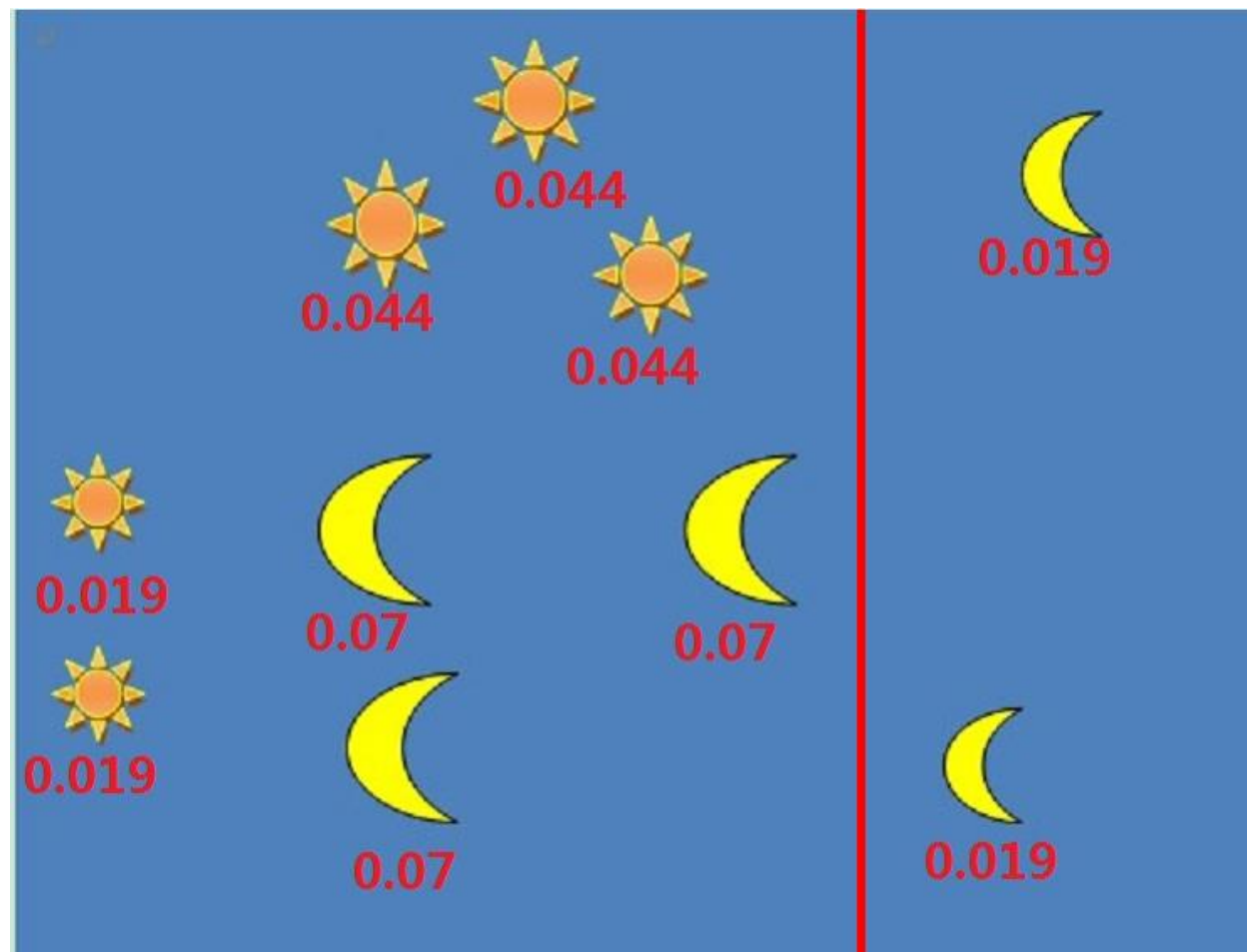
Adaboost:例子

第二轮：M2

更新被正确分类的元组的权重：

$$0.07 * \frac{\text{error}(M_2)}{1 - \text{error}(M_2)} = 0.019$$

$$0.167 * \frac{\text{error}(M_2)}{1 - \text{error}(M_2)} = 0.044$$



Adaboost:例子

归一化:

旧的权重之和:

$$0.167*3+0.07*4=0.99$$

新的权重之和:

$$0.019*4+0.044*3+0.07*3=0.418$$

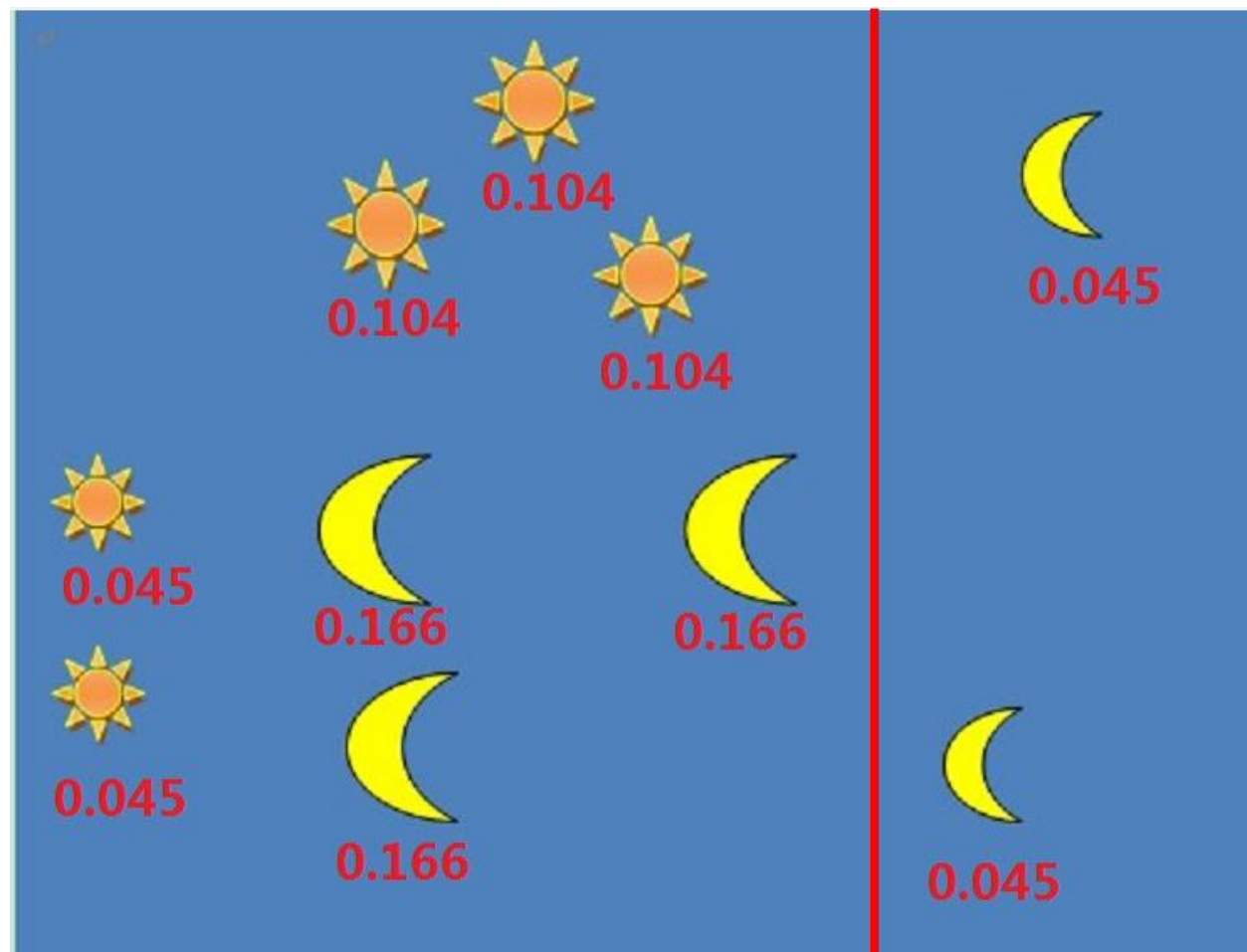
错误地分类:

$$0.07 * \frac{0.99}{0.418} = 0.166$$

正确地分类:

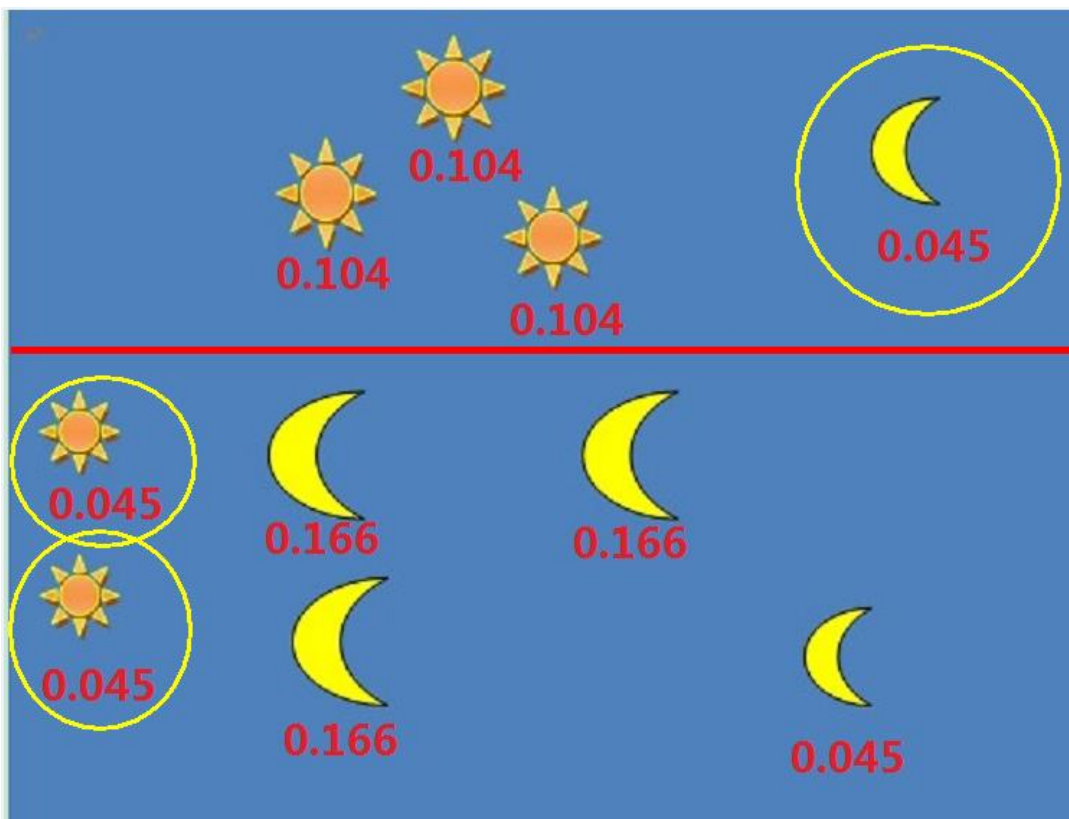
$$0.019 * \frac{0.99}{0.418} = 0.045$$

$$0.044 * \frac{0.99}{0.418} = 0.104$$



Adaboost:例子

第三轮：M3



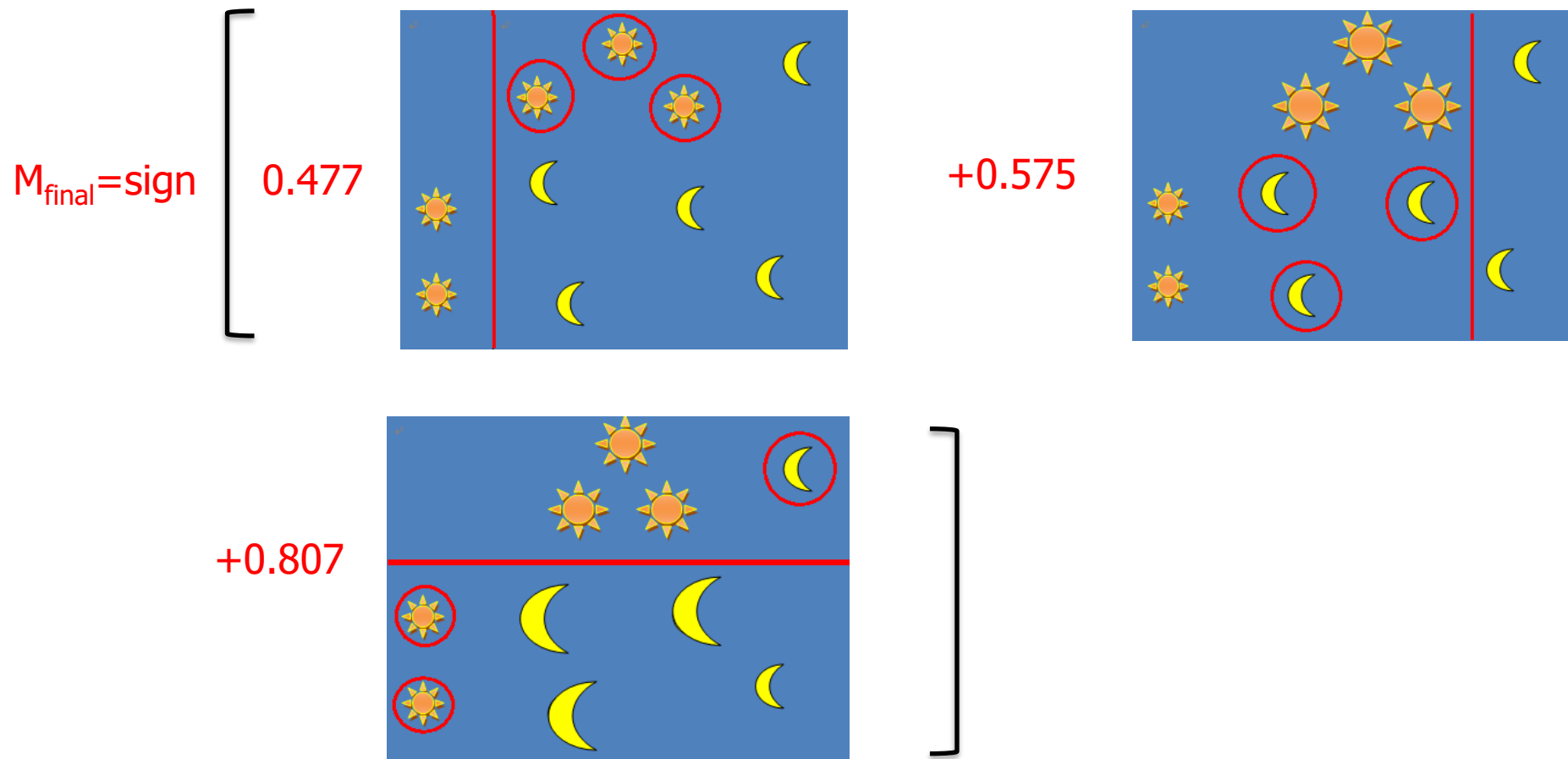
M3的错误率:

$$\begin{aligned} error(M_3) &= \sum_{j=1}^{10} w_j \times err(X_j) \\ &= 0.045 + 0.045 + 0.045 \\ &= 0.135 \end{aligned}$$

M3的表决权重:

$$W_3 = \log \frac{1 - error(M_3)}{error(M_3)} = 0.807$$

Adaboost:例子



梯度提升 (Gradient Boosting) 方法的基本思想

- 可加模型 (Additive Model) :

$$F(\mathbf{x}) = \sum_{j=1}^k f_j(\mathbf{x}) \quad \xrightarrow{\text{加权版本}} \quad F(\mathbf{x}) = \sum_{j=1}^k \rho_j \cdot f_j(\mathbf{x})$$

- 损失函数: $\mathbb{X} = \{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_n, y_n)\}$

$$L(F(\mathbb{X})) = \sum_{i=1}^n l(F(\mathbf{x}_i), y_i) \quad \xrightarrow{\text{目标}} \quad F^* = \arg \min_F L(F(\mathbb{X}))$$
$$= \arg \min_F \sum_{i=1}^n [y_i - F(\mathbf{x}_i)]^2$$

梯度提升机 (Gradient Boosting Machine)

➤ 随机选取初始的基模型 $f_0(\mathbf{x})$, 令 $F_0 = f_0(\mathbf{x})$ $\mathbb{X} = \{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_n, y_n)\}$

➤ for $j=1, 2, \dots, k$

step 1: $\tilde{y}_i = -\frac{\partial L(F(\mathbf{x}_i))}{\partial F} \big|_{F=F_{j-1}}$ pseudo response

$$F(\mathbf{x}) = \sum_{j=1}^k \rho_j \cdot f_j(\mathbf{x})$$

step 2: $\theta^* = \arg \min_{\theta} \sum_{i=1}^n [\tilde{y}_i - f_j(\mathbf{x}_i)]^2$

$$F^* = \arg \min_F L(F(\mathbb{X}))$$

step 3: $\rho_j^* = \arg \min_{\rho_j} \sum_{i=1}^n L[y_i, F_{j-1}(\mathbf{x}_i) + \rho_j \cdot f_j(\mathbf{x}_i)]^2$

$$= \arg \min_F \sum_{i=1}^n [y_i - F(\mathbf{x}_i)]^2$$

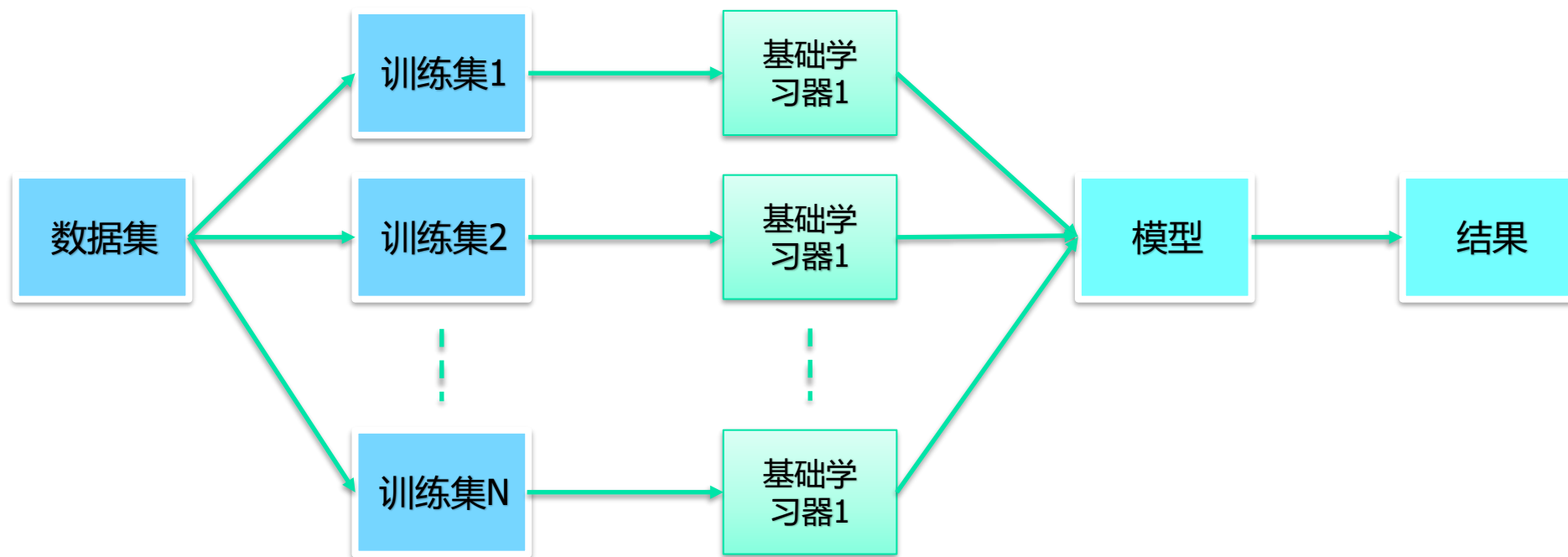
step 4: $F_j = F_{j-1} + \rho_j^* f_j$

XGBoost & LightGBM

- Tianqi Chen, Carlos Guestrin(2016). XGBoost: A Scalable Tree Boosting System. KDD'16.
- Guolin Ke, Qi Meng, Thomas Finley, Taifeng Wang, Wei Chen, Weidong Ma, Qiwei Ye, Tie-Yan Liu. "LightGBM: A Highly Efficient Gradient Boosting Decision Tree". Advances in Neural Information Processing Systems 30 (NIPS 2017), pp. 3149-3157.

堆叠法 (stacking) 基本原理

- 异构集成学习方法
 - 组成堆叠法的基础学习器一般互不相同



堆叠法的关键问题

- 保证不同个体学习数据的独立性
 - 对数据集进行随机划分
- 保证不同个体的多样性
 - 对不同的数据集采取不同的模型训练
- 集成策略的选择
 - 将多个模型的输出作为一个模型的输入

致谢

- 部分图表、文字来自教材、互联网等，仅供公益性的学习参考，在此表示感谢！如有版权要求请联系：yym@hit.edu.cn，谢谢！