

大数据分析

Scalable Machine Learning
decision tree

刘盛华

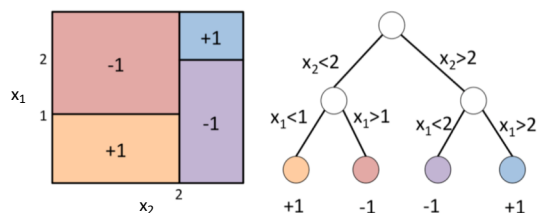
Outline

- 决策树 (Decision Tree)
- 随机森林 (Random Forest)
- 梯度提升树 (Gradient Boosted Decision Tree (GBDT))

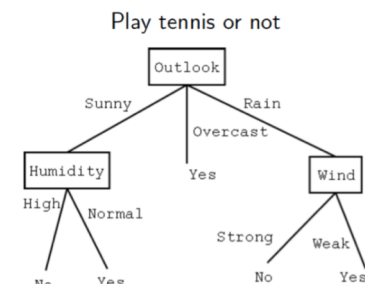
Cho-Jui Hsieh, UC Davis ECS171: Machine Learning

Decision Tree

- 每个节点负责检查特征 x_i :
 - 若 $x_i < \text{threshold}$, 选择左子树
 - 若 $x_i \geq \text{threshold}$, 选择右子树



例子



Decision Tree

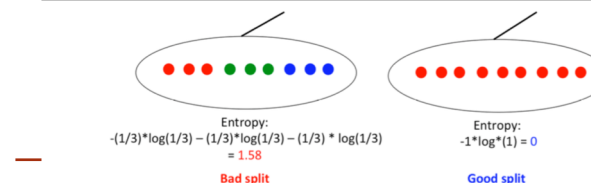
- 计算能力:
 - 决策树是**非线性**分类器
 - 决策树更具备**可解释性**
 - 决策树天然适合处理**类别型特征**
- 计算效率:
 - 训练: **较慢**
 - 预测: **较快**
 - 决策树有 h 步操作operations (h : 树的深度, 通常 ≤ 15)

划分节点

- 分类树: 依照使数据的熵最大的准则划分节点
- S 表示一个节点中的全部样本, 每个样本的标签 $c = 1, \dots, C$:

$$\text{Entropy} : H(S) = - \sum_{c=1}^C p(c) \log p(c),$$

- $p(c)$ 表示所有样本中属于类别 c 的样本所占比例.
 - 若所有样本属于同一类别, Entropy=0
 - 若 $p(1) = \dots = p(C)$, Entropy最大



信息增益

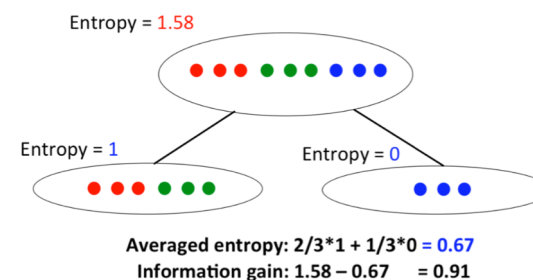
- 划分 $S \rightarrow S_1, S_2$ 的熵的平均值为:

$$\frac{|S_1|}{|S|} H(S_1) + \frac{|S_2|}{|S|} H(S_2)$$

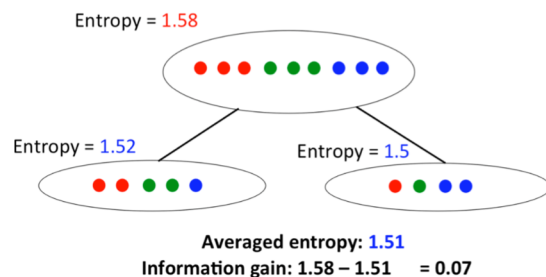
- 信息增益: 评价划分方式的好坏

$$H(S) - \left(\frac{|S_1|}{|S|} H(S_1) + \frac{|S_2|}{|S|} H(S_2) \right)$$

信息增益



信息增益

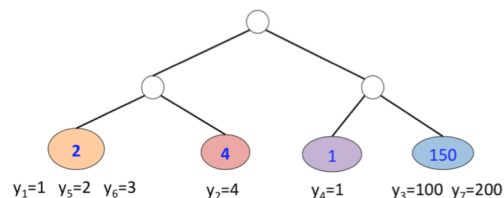


节点划分

- 给定当前节点，如何选择**最优划分**方法？
- 对所有特征和阈值：
 - 计算按照每个特征及阈值划分后的信息增益
 - 选择最优划分 (**最大化信息增益**)
- 对 n 个样本和 d 个特征: 时间复杂度 $O(nd)$

回归树

- 为每个叶子节点赋一个数值
- 通常每个叶子节点取值为它包含的所有 y 个样本的平均 (最小化平方误差)



回归树

Objective function:

$$\min_F \frac{1}{n} \sum_{i=1}^n (y_i - F(x_i))^2 + (\text{Regularization})$$

The quality of partition $S = S_1 \cup S_2$ can be computed by the objective function:

$$\sum_{i \in S_1} (y_i - y^{(1)})^2 + \sum_{i \in S_2} (y_i - y^{(2)})^2,$$

where $y^{(1)} = \frac{1}{|S_1|} \sum_{i \in S_1} y_i$, $y^{(2)} = \frac{1}{|S_2|} \sum_{i \in S_2} y_i$ 每个划分集的方差最小

Find the best split:

Try all the features & thresholds and find the one with **minimal objective function**

决策树的超参数

- 最大树深: (通常 ~ 10)
- 分割节点所需的最小样本数: (10, 50, 100)
- 单决策树通常效果有限...
- 是否可以构建多棵决策树, 然后集成?

Outline

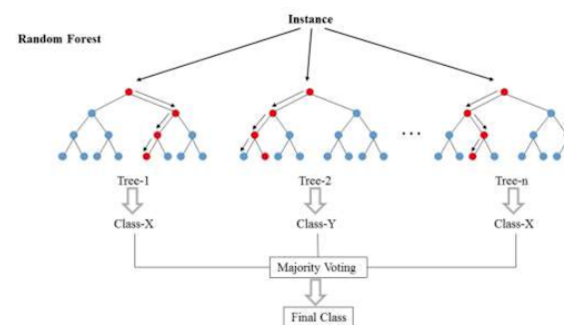
- 决策树 (Decision Tree)
- 随机森林 (Random Forest)
- 梯度提升树 (Gradient Boosted Decision Tree (GBDT))

Cho-Jui Hsieh, UC Davis ECS171: Machine Learning

随机森林

- 随机森林 (多决策树通过自助法集成 (Bootstrap ensemble)):
 - 创建 T 棵树
 - 对每棵树, 采样一个数据子集 S_i 和特征子集 D_i 用来训练
 - 预测: 取 T 棵树的结果的平均
- 随机森林的优点:
 - 避免过拟合
 - 改善稳定性和准确率
- 有很好的软件包:
 - R: "randomForest" package
 - Python: sklearn

例子



通过MapReduce构建决策树

- Parallel Learner for Assembling Numerous Ensemble Trees [Panda et al., VLDB '09]
 - A sequence of MapReduce jobs that builds a decision tree
 - Spark MLlib Decision Trees are based on PLANET

Outline

- 决策树 (Decision Tree)
- 随机森林 (Random Forest)
- 梯度提升树 (Gradient Boosted Decision Tree (GBDT))

Cho-Jui Hsieh, UC Davis ECS171: Machine Learning

提升树 (Boosted Decision Tree)

- Minimize loss $\ell(y, F(x))$ with $F(\cdot)$ being ensemble trees

$$F^* = \operatorname{argmin}_F \sum_{i=1}^n \ell(y_i, F(x_i)) \quad \text{with} \quad F(x) = \sum_{m=1}^T f_m(x)$$

(each f_m is a decision tree)

- Direct loss minimization: at each stage m , find the best function to minimize loss

- solve $f_m = \operatorname{argmin}_{f_m} \sum_{i=1}^N \ell(y_i, F_{m-1}(x_i) + f_m(x_i))$
- update $F_m \leftarrow F_{m-1} + f_m$

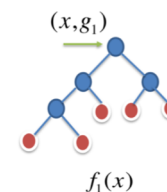
$F_m(x) = \sum_{j=1}^m f_j(x)$ is the prediction of x after m iterations.

- Two problems:
 - Hard to implement for general loss
 - Tend to overfit training data

梯度提升树 (GBDT)

Key idea:

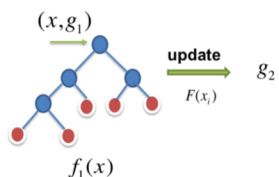
- Each base learner is a decision tree
- Each regression tree approximates the functional gradient $\frac{\partial \ell}{\partial F}$



梯度提升树 (GBDT)

Key idea:

- Each base learner is a decision tree
- Each regression tree approximates the functional gradient $\frac{\partial \ell}{\partial F}$

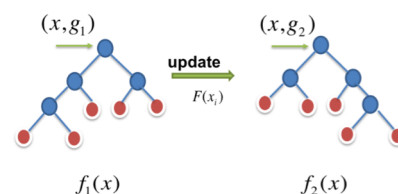


$$F_{m-1}(x_i) = \sum_{j=1}^{m-1} f_j(x_i) \quad g_m(x_i) = \left. \frac{\partial \ell(y_i, F(x_i))}{\partial F(x_i)} \right|_{F(x_i)=F_{m-1}(x_i)}$$

梯度提升树 (GBDT)

Key idea:

- Each base learner is a decision tree
- Each regression tree approximates the functional gradient $\frac{\partial \ell}{\partial F}$

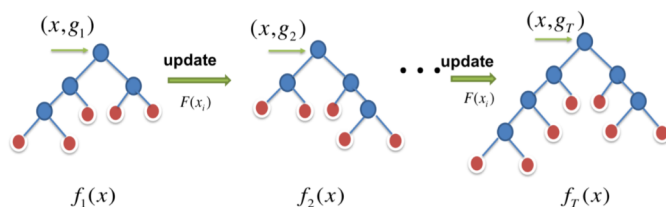


$$F_{m-1}(x_i) = \sum_{j=1}^{m-1} f_j(x_i) \quad g_m(x_i) = \left. \frac{\partial \ell(y_i, F(x_i))}{\partial F(x_i)} \right|_{F(x_i)=F_{m-1}(x_i)}$$

梯度提升树 (GBDT)

• Key idea:

- Each base learner is a decision tree
- Each regression tree approximates the functional gradient $\frac{\partial \ell}{\partial F}$

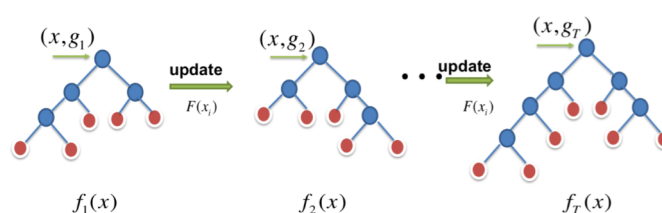


$$F_{m-1}(x_i) = \sum_{j=1}^{m-1} f_j(x_i) \quad g_m(x_i) = \left. \frac{\partial \ell(y_i, F(x_i))}{\partial F(x_i)} \right|_{F(x_i)=F_{m-1}(x_i)}$$

梯度提升树 (GBDT)

• Key idea:

- Each base learner is a decision tree
- Each regression tree approximates the functional gradient $\frac{\partial \ell}{\partial F}$



Final prediction $F(x_i) = \sum_{j=1}^T f_j(x_i)$

Questions?