

分类与回归树



赵海涛

haitaozhao@ecust.edu.cn



Josh Gordon

random-forests

Unfollow

Git is complicated ͇(ツ)͇

📍 NYC

🔗 twitter.com/random_forests

Block or report user

Organizations



CART

Training Data

Color	Diameter	Label
Green	3	Apple
Yellow	3	Apple
Red	1	Grape
Red	1	Grape
Yellow	3	Lemon

Training Data

Color	Diameter	Label
Green	3	Apple
Yellow	3	Apple
Red	1	Grape
Red	1	Grape
Yellow	3	Lemon

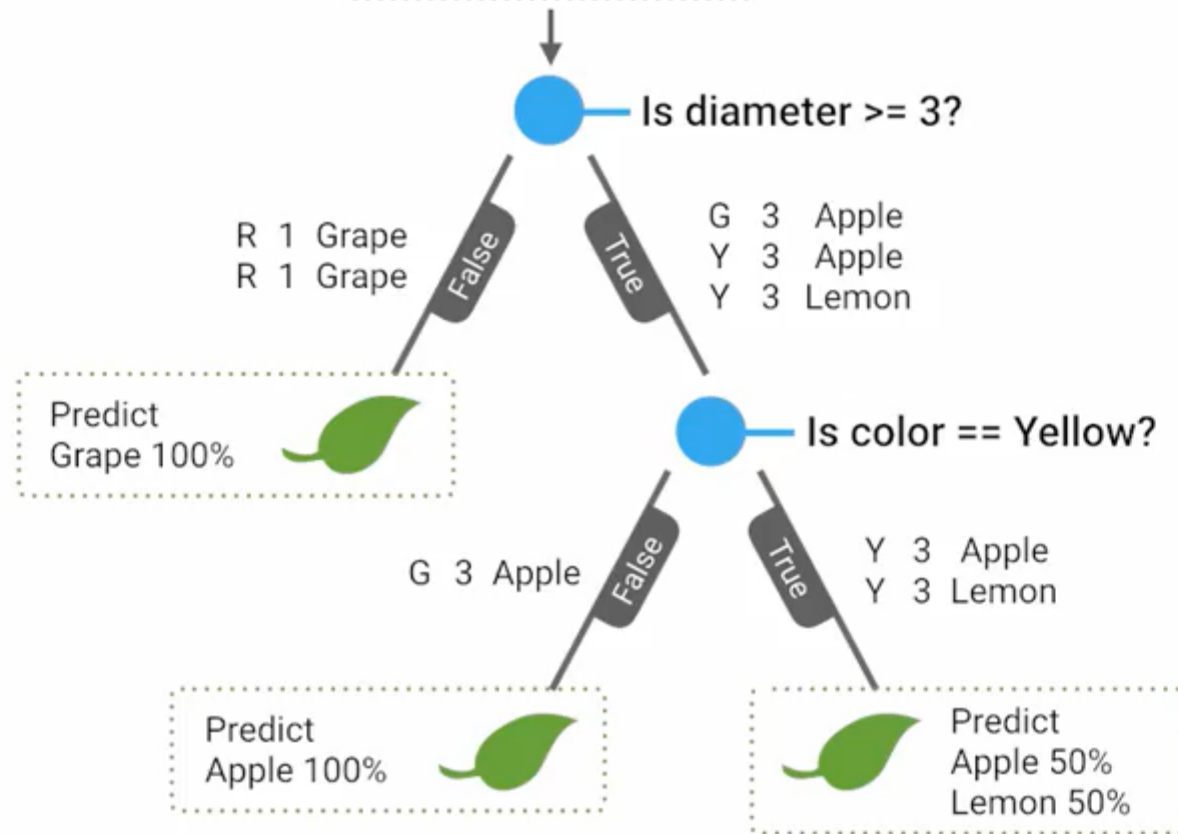
A categorical attribute

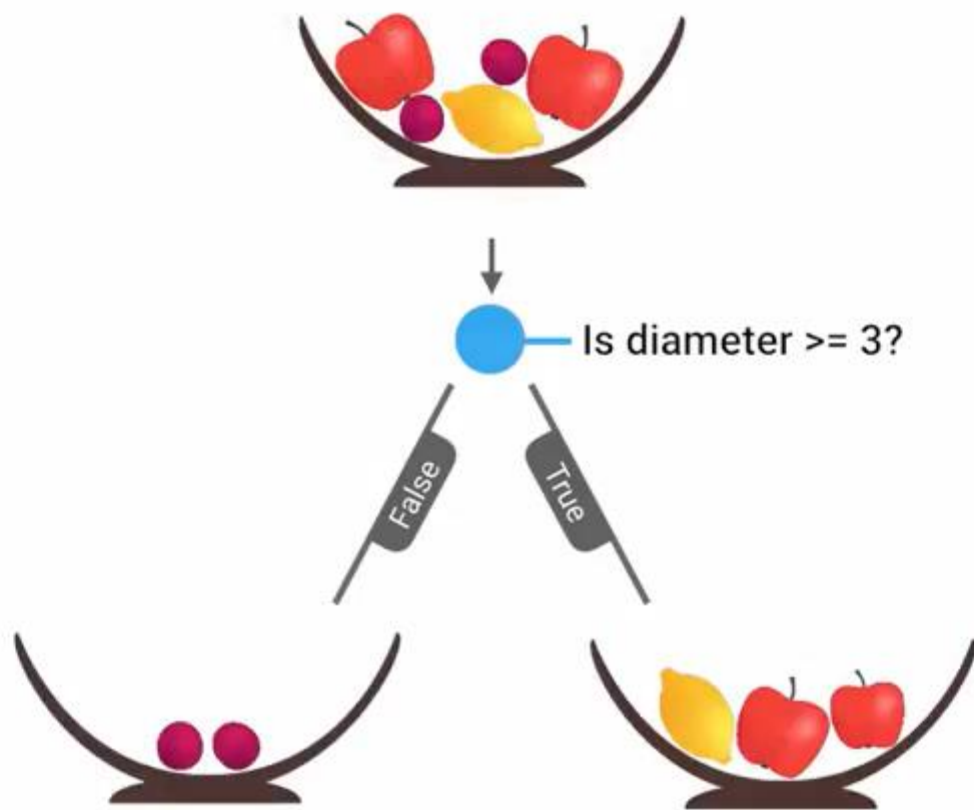
Training Data

Color	Diameter	Label
Green	3	Apple
Yellow	3	Apple
Red	1	Grape
Red	1	Grape
Yellow	3	Lemon

A numeric attribute

Color	Diam	Label
Green	3	Apple
Yellow	3	Apple
Red	1	Grape
Red	1	Grape
Yellow	3	Lemon





Possible questions

Is the color green?

Is the diameter ≥ 3 ?

Is the color yellow?

TRUE

Green 3 Apple

FALSE

Yellow 3 Apple
Yellow 3 Lemon

Green	3	Apple
Yellow	3	Apple
Yellow	3	Lemon

Possible questions

Is the color green?

Is the diameter ≥ 3 ?

Is the color yellow?

TRUE

FALSE

信息增益(Information Gain)

- 定义5.2 (信息增益):特征 A 对训练数据集 D 的信息增益, $g(D, A)$, 定义为集合 D 的经验熵 $H(D)$ 与特征 A 给定条件下 D 的经验条件熵 $H(D|A)$ 之差, 即

$$g(D, A) = H(D) - H(D|A)$$

- (Information gain)表示得知特征 X 的信息而使得类 Y 的信息的不确定性减少的程度。
- 一般地, 熵 $H(Y)$ 与条件熵 $H(Y|X)$ 之差称为互信息 (mutual information)
- 决策树学习中的信息增益等价于训练数据集中类与特征的互信息。

信息增益(Information Gain)

- 输入：训练数据集 D 和特征 A ;
- 输出：特征 A 对训练数据集 D 的信息增益 $g(D, A)$

1、计算数据集 D 的经验熵 $H(D)$

$$H(D) = - \sum_{k=1}^K \frac{|C_k|}{|D|} \log_2 \frac{|C_k|}{|D|}$$

2、计算特征 A 对数据集 D 的经验条件熵 $H(D|A)$

$$H(D|A) = \sum_{i=1}^n \frac{|D_i|}{|D|} H(D_i) = - \sum_{i=1}^n \frac{|D_i|}{|D|} \sum_{k=1}^K \frac{|D_{ik}|}{|D_i|} \log_2 \frac{|D_{ik}|}{|D_i|}$$

3、计算信息增益

$$g(D, A) = H(D) - H(D|A)$$

信息增益(Information Gain)

- 输入：父节点数据集 D_p 和对其的划分得到的两个子数据集 D_{left} , D_{right} ;
- 输出：信息增益 $IG(D_p)$

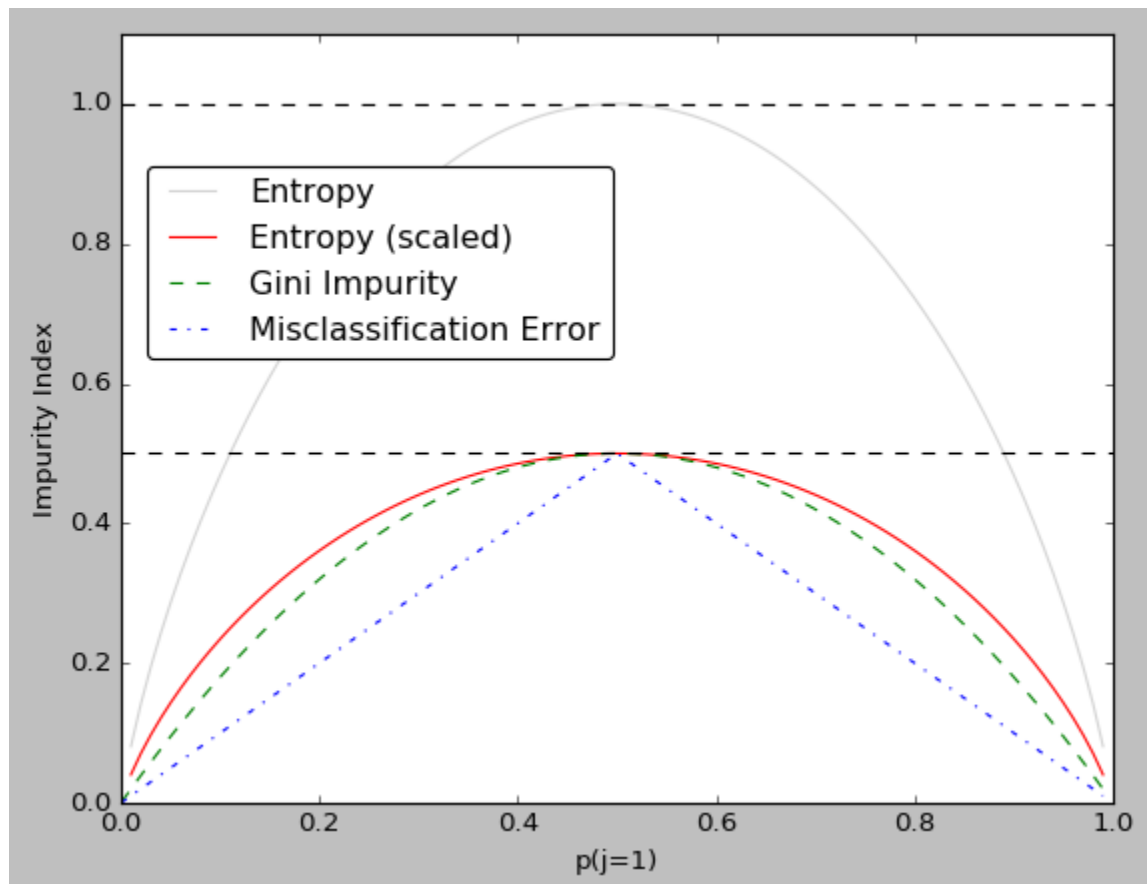
计算信息增益

$$IG(D_p) = I(D_p) - \frac{N_{left}}{N_p} I(D_{left}) - \frac{N_{right}}{N_p} I(D_{right})$$

- 可能的不纯度指数
 - ① Entropy: $\sum_{j=1}^K p_j \ln \frac{1}{p_j}$.
 - ② Misclassification rate: $1 - \max_j p_j$.
 - ③ Gini index: $\sum_{j=1}^K p_j(1-p_j) = 1 - \sum_{j=1}^K p_j^2$.

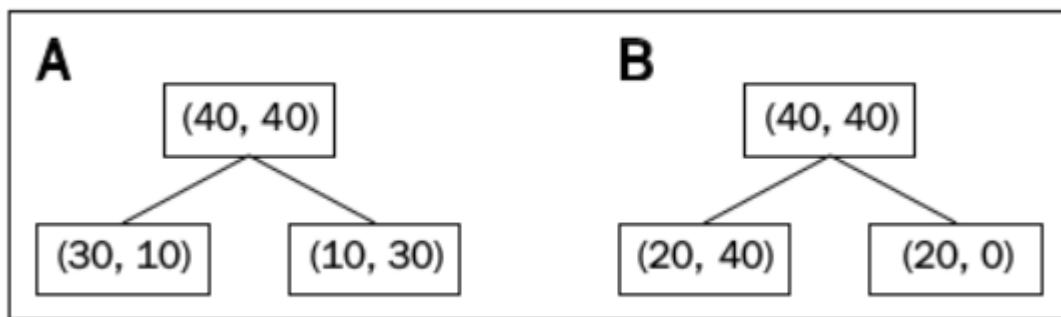
不纯度指数(Impurity Index)

- ① Entropy: $\sum_{j=1}^K p_j \ln \frac{1}{p_j}$.
- ② Misclassification rate: $1 - \max_j p_j$.
- ③ Gini index: $\sum_{j=1}^K p_j(1-p_j) = 1 - \sum_{j=1}^K p_j^2$.



不纯度指数(Impurity Index)

- ① Entropy: $\sum_{j=1}^K p_j \ln \frac{1}{p_j}$.
- ② Misclassification rate: $1 - \max_j p_j$.
- ③ Gini index: $\sum_{j=1}^K p_j(1-p_j) = 1 - \sum_{j=1}^K p_j^2$.



Color	Diam	Label
Green	3	Apple
Yellow	3	Apple
Red	1	Grape
Red	1	Grape
Yellow	3	Lemon

Impurity = 0.64

Impurity = 0.62

Yellow	3	Apple
Red	1	Grape
Red	1	Grape
Yellow	3	Lemon

Is color green?

False

Green 3 Apple

True

Impurity = 0

Information Gain = $0.64 - 0.5 = 0.14$

Color	Diam	Label
Green	3	Apple
Yellow	3	Apple
Red	1	Grape
Red	1	Grape
Yellow	3	Lemon

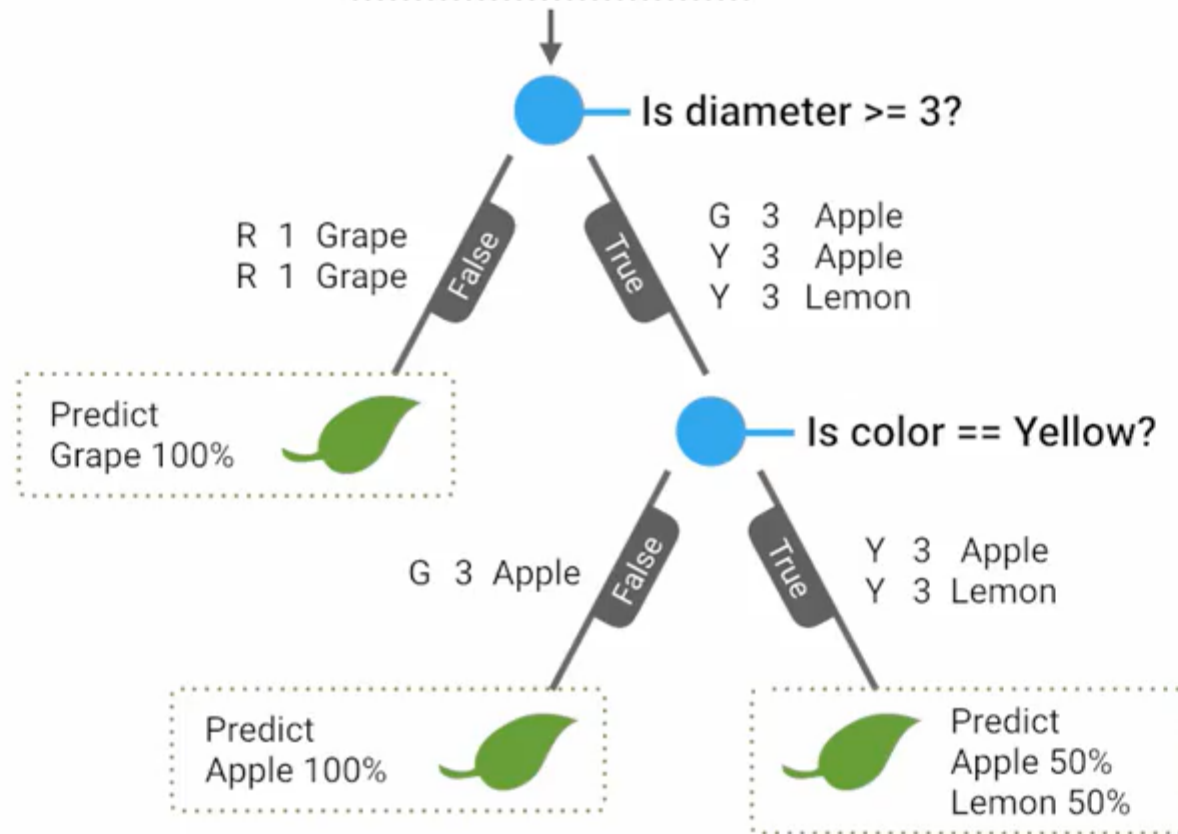


Information Gain

Question	Gain
Color == Green?	0.14
Diameter >= 3?	0.37
Color == Yellow?	0.17
Color == Red?	0.37
Diameter >= 1?	0

Here there is a tie. We'll pick the first.

Color	Diam	Label
Green	3	Apple
Yellow	3	Apple
Red	1	Grape
Red	1	Grape
Yellow	3	Lemon



谢谢各位同学！