

decision tree

ID3 (Iterative Dichotomiser 3) :

- 使用信息增益作为分割数据的标准。
- 能够处理分类数据。

C4.5:

- 是ID3的后继者，使用信息增益比率来选择特征。
- 能够处理分类和连续数据。
- 对缺失数据有处理能力。

CART (Classification and Regression Trees) :

- 可以用于分类和回归（名字中的R代表回归）。
- 使用基尼不纯度（Gini impurity）或均方误差（Mean Squared Error）作为分割数据的标准。
- 生成的是二叉树。

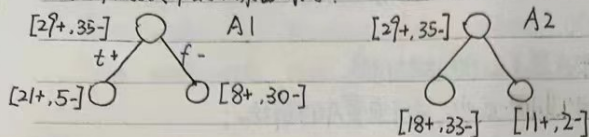
信息熵

$$H(A) = -P(A) \cdot \log_2 P(A)$$

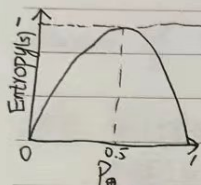
信息增益就是信息熵的差值

Top-Down Induction of Decision Trees..

建立最有效率的决策数树.



样本熵. Sample Entropy



S : 训练样本

p_0 : positive examples 在 S 中占比

p_0 : negative examples 在 S 中占比

样本 S 的不纯度熵 $H(S) = -p_0 \log_2 p_0 - p_0 \log_2 p_0$

Entropy $H(X)$ of a random variable X

$$H(X) = -\sum_{i=1}^n P(X=i) \log_2 P(X=i)$$

↓ 期望 bits 数去 encode 随机取出的 X 值

信息理论: 大多数效率代码需要 $-\log_2 P(X=i)$ bits 去 ^编 码信息 $X=i$.

则期望值: $\sum_{i=1}^n P(X=i) (-\log_2 P(X=i))$

具体条件下熵 $H(X|Y=v)$ 的 X with given $Y=v$

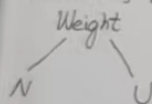
$$H(X|Y=v) = -\sum_{i=1}^n P(X=i|Y=v) \log_2 P(X=i|Y=v)$$

The following data set will be used to learn a decision tree for predicting whether students are lazy (L) or diligent (D) based on their weight (Normal or Underweight), their eye color (Amber or Violet) and the number of eyes they have (2 or 3 or 4).

| Weight | Eye Color | Num. Eyes | Output |
|--------|-----------|-----------|--------|
| N | A | 2 | L |
| N | V | 2 | L |
| N | V | 2 | L |
| U | V | 3 | L |
| U | V | 3 | L |
| U | A | 4 | D |
| N | A | 4 | D |
| N | V | 4 | D |
| U | A | 3 | D |
| U | A | 3 | D |

The following numbers may be helpful as you answer this problem without using a calculator:
 $\log_2 0.1 = -3.32$, $\log_2 0.2 = -2.32$, $\log_2 0.3 = -1.73$, $\log_2 0.4 = -1.32$, $\log_2 0.5 = -1$.
 *You don't need to show the derivation for your answers in this problem.

2. Weight:
 $S: [5, 5-] \quad E_1 = 1$



$$[3, 2-] \quad [2, 3-] \\ E_2 = \frac{2}{5} \log_2 \frac{2}{5} + \frac{3}{5} \log_2 \frac{3}{5} = 0.992 \\ = 0.438 + 0.464 = 0.902$$

$$\text{Gain}(S, \text{Weight}) = E_1 - \frac{1}{2} E_2 - \frac{1}{2} E_3 = 1 - \frac{1}{2} (0.992) - \frac{1}{2} (0.992) = 0.504$$

1. (3 pts) What is the conditional entropy $H(\text{Eye Color} | \text{Weight} = N)$?

$$H(\text{Eye Color} | \text{Weight} = N) = -P(\text{Eye} = A | \text{Weight} = N) \cdot \log_2 P(\text{Eye} = A | \text{Weight} = N) \\ - P(\text{Eye} = V | \text{Weight} = N) \cdot \log_2 P(\text{Eye} = V | \text{Weight} = N) = -0.4 \cdot \log_2 0.4 - 0.6 \cdot \log_2 0.6 = 0.966$$

2. (3 pts) What attribute would the ID3 algorithm choose to use for the root of the tree (pruning)? Eye Color

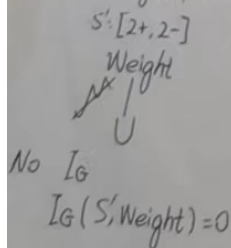
$$E_1 = 1, E_2 = \frac{1}{5} \log_2 \frac{1}{5} + \frac{4}{5} \log_2 \frac{4}{5} = 0.722, E_3 = \frac{2}{5} \log_2 \frac{2}{5} + \frac{3}{5} \log_2 \frac{3}{5} = 0.992 \\ \text{Gain}(S, \text{Eye Color}) = E_1 - \frac{1}{5} E_2 - \frac{4}{5} E_3 = 1 - \frac{1}{5} (0.722) - \frac{4}{5} (0.992) = 0.278$$

3. (4 pts) Draw the full decision tree learned for this data (no pruning).

$$\therefore \text{Gain}(S, \text{Num Eyes}) \text{ is the biggest. } \therefore \text{Num of Eyes should be the root of the tree.} \\ \text{Gain}(S, \text{Num Eyes}) = E_1 - \frac{3}{10} E_2 - \frac{4}{10} E_3 - \frac{3}{10} E_4 = 1 - \frac{3}{10} (0.992) - \frac{4}{10} (0.992) - \frac{3}{10} (0.992) = 0.6$$

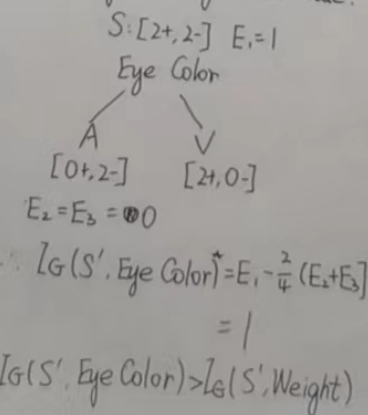
3. \therefore the output of $\text{Num Eyes} = 2$ is L, the output of $\text{Num Eyes} = 4$ is D.
 \therefore Let's talk about the condition that $\text{Num Eyes} = 3$.

① Use Weight for next Node



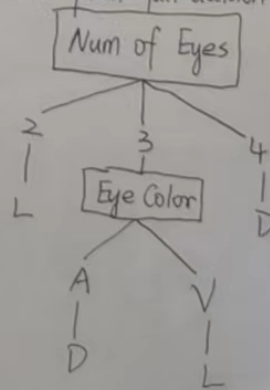
$$IG(S', \text{Weight}) = 0$$

② Use Eye Color for next Node.



$$IG(S', \text{Eye Color}) = E_1 - \frac{2}{4} (E_2 + E_3) = 1 - \frac{2}{4} (0.992 + 0.992) = 0.504 \\ IG(S', \text{Eye Color}) > IG(S', \text{Weight})$$

\therefore the final full decision tree:



Weight

$S: [5+, 5-]$ $E_1 = 1$

Weight

$\begin{array}{c} 5 \quad 5 \\ \swarrow \quad \searrow \\ N \quad U \end{array}$

$[3+, 2-]$ $[2+, 3-]$

$$E_2 = \frac{3}{5} \log_2 \frac{3}{5} - \frac{2}{5} \log_2 \frac{2}{5} = 0.966$$

$$E_3 = E_2 = 0.966$$

$$\text{Gain}(S, \text{Weight}) = E_1 - \frac{1}{2}E_2 - \frac{1}{2}E_3 = 0.034$$

3. $S': [2+, 2-]$ $E_1 = 1$

Weight

$\begin{array}{c} 4 \\ | \\ U \end{array}$

$[2+, 2-]$

$$E_2 = 1$$

$$\text{Gain} = E_1 - \frac{1}{2}E_2 = 0$$

$S': [2+, 2-]$ $E_1 = 1$

Eye color

$\begin{array}{c} 2 \quad 2 \\ \swarrow \quad \searrow \\ A \quad V \end{array}$

$[10+, 2-]$ $[2+, 10-]$

$$E_2 = 0$$

$$E_3 = 0$$

$$\text{Gain} = E_1 - \frac{1}{2}E_2 - \frac{1}{2}E_3 = 1$$

$$I_G(S', \text{Eye Color}) > I_G(S', \text{Weight})$$