

der-23 04 | 03 | 21

→ statistical measures (contd.)
→ construction of DT

② Gain ratio:

- C4.5, a successor of ID3, uses an extension to info. gain known as gain ratio.
- It overcomes the bias factor encountered in info. gain.
- It applies a kind of normalization to info gain using split informations.

$$\text{Gain Ratio} = \frac{\text{Gain}(A)}{\text{Split info}(A)}$$

$$\text{Split info}_A(D) = \sum_{j=1}^v \frac{|D_j|}{|D|} \log_2 \left(\frac{|D_j|}{|D|} \right)$$

This value represents the potential information generated by splitting the training data set D , into v partitions, corresponding to v outcomes of the test attribute A .

* The attribute with the maximum gain ratio is selected as the splitting attribute.

Gain (Age): 0.246 bits

Gain (Income) 0.029 bits

Gain (student) 0.151 bits

Gain (CR) = 0.048 "

split info (income) :

↓
low (0.4)
med (0.6)
high (0.4)

$$\begin{aligned} \text{split info (income)} &= -\frac{4}{14} \log_2 \frac{4}{14} - \frac{6}{14} \log_2 \frac{6}{14} - \frac{4}{14} \log_2 \frac{4}{14} \\ &= \underline{\underline{1.556}} \end{aligned}$$

$$\text{Gain ratio (income)} = \frac{0.029}{1.556} = \underline{\underline{0.0188}}$$

③ Gini Index:

→ This measure is used in CART.

→ It considers a binary split for each attribute.

→ The Gini Index measures the impurity of D

$$\boxed{\text{Gini}(D) = 1 - \sum_{i=1}^m P_i^2}$$

where P_i is the probability that a tuple in D belongs to class C_i as is estimated by $\hat{P}_{(i,D)}$

$$\boxed{\text{Gini}_A(D) = \frac{|D_1|}{|D|} \text{Gini}(D_1) + \frac{|D_2|}{|D|} \text{Gini}(D_2)}$$

$$\boxed{\Delta \text{Gini}(A) = \text{Gini}(D) - \text{Gini}_A(D)}$$

The reduction in impurity that would be incurred by a binary split on a discrete or continuous valued attribute A .

→ The attribute that maximizes the redn. in impurity is selected as the splitting attribute.

$$\text{Prob: } \text{Gini}(D) = 1 - \left(\frac{9}{14}\right)^2 - \left(\frac{5}{14}\right)^2$$

buy car
yes buy car
no
 (0.9) (0.4)

$$= 0.459$$

=====

Gini (Income)

Income (low med high)
 (0.4) (0.6) (0.4)

$$\Rightarrow \begin{array}{ll} \{ \text{low, med} \} & \{ \text{high} \} \\ \{ \text{low} \} & \{ \text{high} \} \\ \{ \text{med, high} \} & \{ \text{low} \} \\ \{ \text{high, low} \} & \{ \text{medium} \} \end{array} - \text{cat(1)} \quad \text{cat(2)} \quad \text{cat(3)}$$

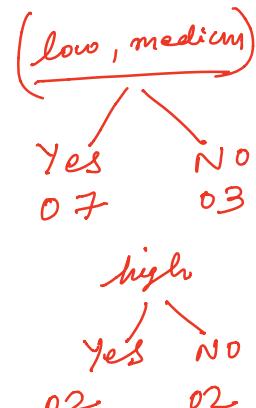
$\text{Gini} \left\{ \frac{\{ \text{low, med} \}}{10}, \frac{\{ \text{high} \}}{4} \right\}$

$$= \frac{10}{14} \underline{\text{Gini}(D_1)} + \frac{4}{14} \underline{\text{Gini}(D_2)}$$

$$= \frac{10}{14} \left(1 - \left(\frac{7}{10} \right)^2 - \left(\frac{3}{10} \right)^2 \right) +$$

$$\frac{4}{14} \left(1 - \left(\frac{2}{4} \right)^2 - \left(\frac{2}{4} \right)^2 \right)$$

$$= \underline{0.443}$$



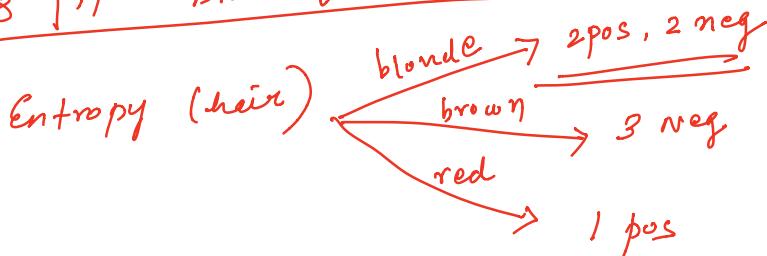
$$\text{Gini}_{\text{Cat}(2)} \{ \text{low, high} \} \{ \text{med} \} = \underline{0.458}$$

$$\text{Gini}_{\text{Cat}(3)} \{ \text{med, high} \} \{ \text{low} \} = \underline{0.45}.$$

As Gini is to be minimized, therefore category (1) has the lowest Gini, which is the preferable binary split.

Problem to construct DT. considering Entropy as the statistical measure.

	Name	Hair	Height	Weight	Letters	Sunburn(res)
1	A	Blonde ✓	avg	light	NO	pos
2	B	Blonde ✓	tall	avg	yes	neg
3	C	Brown	short	avg	yes	neg
4	D	Blonde ✓	short	avg.	NO	pos
5	E	Red	avg	heavy	NO	pos
6	F	Brown	tall	heavy	NO	neg
7	G	Brown	avg	heavy	NO	neg
8	H	Blonde ✓	short	light	yes	neg



$$\begin{aligned}
 \text{Entropy} : & \frac{4}{8} \left(-\frac{2}{4} \log_2 \frac{2}{2} - \frac{2}{4} \log_2 \frac{2}{2} \right) + \\
 & \frac{3}{8} \left(-0 - \frac{3}{3} \log_2 \frac{3}{3} \right) + \\
 & \frac{1}{8} \left(-\frac{1}{1} \log_2 \frac{1}{1} - 0 \right) \\
 = & \underline{\underline{0.5}}
 \end{aligned}$$

Entropy (hair):	0.51
Ent (height):	0.69
Ent (weight):	0.94
Ent (lotion):	0.61

hair color, blonde — ? ✓
 hair color, red — ?
 hair color, brown — ?

$\langle \text{Blonde}, \text{height} \rangle$

$E(\text{Blonde}, \text{height}):$

tall	avg	short
($\frac{1}{2}$)	($\frac{1}{2}$)	($\frac{1}{2}$)
(neg)	(pos)	$\begin{array}{l} \text{pos}(\frac{1}{2}) \\ \text{neg}(\frac{1}{2}) \end{array}$

$$\frac{1}{4} \left(-\frac{1}{2} \log_2 \frac{1}{2} - 0 \right) + \frac{1}{4} \left(-\frac{1}{2} \log_2 \frac{1}{2} - 0 \right) +$$

$$\frac{2}{4} \left[-\frac{1}{2} \log_2 \frac{1}{2} - \frac{1}{2} \log_2 \frac{1}{2} \right]$$

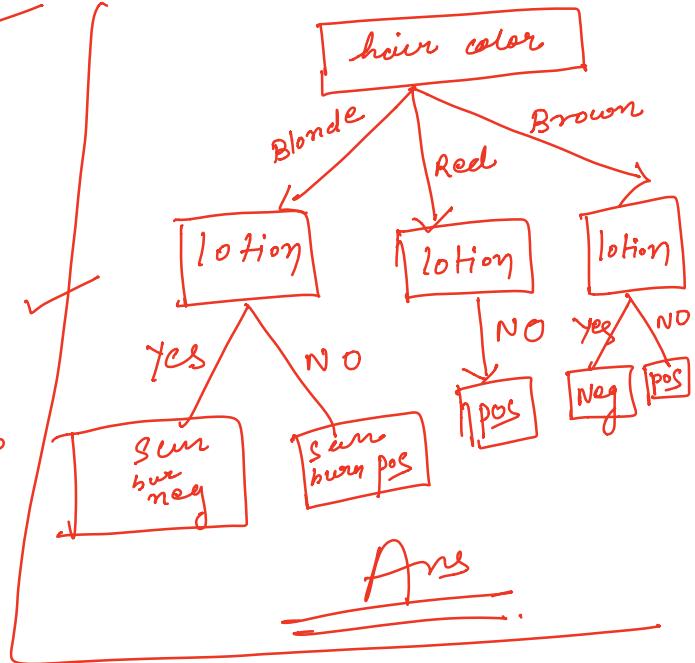
$$= \underline{\underline{0.54}}$$

Entropy (Blonde, wt)

$$= \frac{2}{4} \left(-\frac{1}{2} \log_2 \frac{1}{2} - \frac{1}{2} \log_2 \frac{1}{2} \right) + \frac{2}{4} \left(-\frac{1}{2} \log_2 \frac{1}{2} - \frac{1}{2} \log_2 \frac{1}{2} \right)$$

$$+ \frac{0}{4} \left[-0 - 0 \right]$$

$$= 1$$



Entropy (blonde, lotion)

$$E = \frac{2}{3} \left(-\frac{1}{2} \log_2 \frac{1}{2} - 0 \right) + \frac{2}{3} \left(-0 - \frac{2}{2} \log_2 \frac{1}{2} \right)$$

$$= 0$$

Entropy (Blonde, ht) = 0.5

Entropy (Blonde, wt) = 1

Entropy (Blonde, lotion) = 0

