

Dec-22 02/03/21

Decision Trees

→ Introduction

→ Statistical measures

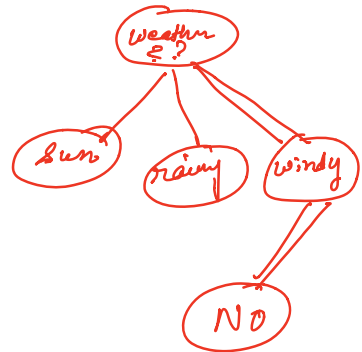
Suppose the question is:

"What is the weather outlook today?" for attribute seln.

↓
"Is it windy?"

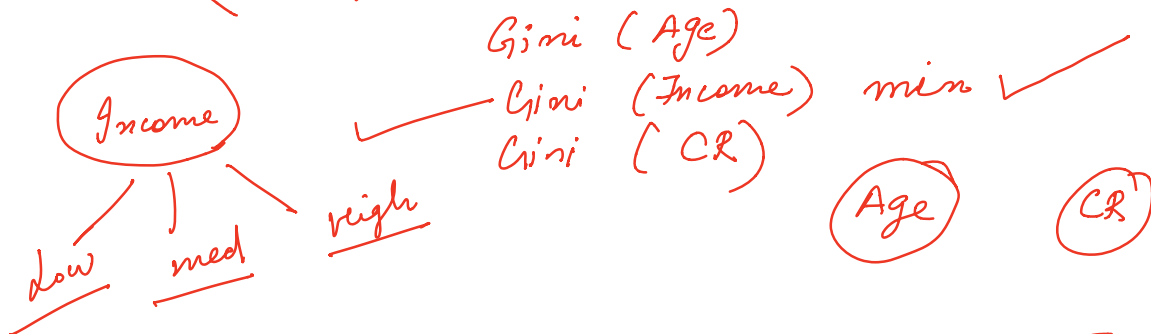
↓
"Should we play golf today?"

↓
"No"



< Age, Income, credit-rating >

< Entropy Info. gain split into Quinlan >



→ DT are otherwise known as Classifications Trees or hierarchical classifiers. These are rule-based learning - top down induction method.

→ Quinlan [1986] introduced ID3 [Interactive Dichotomizer third series] for DT

ID3: Entropy is used to measure the information content of an attribute.

→ C4.5 is the advanced version of ID3

Information gain is used in C4.5.

① Entropy / Information Gain :

→ Entropy is a measure of how much uncertainty is present in the information.

→ Information gain is a measure of how much information the answer to a specific question provides.

Prob Entropy:
$$H(X) = - \sum_{i=1}^n P_i \log_2 P_i$$

where P_i is the probability of occurrence of i

Q Let us consider 3 possible events A, B and C having equal prob. of occurrence. Find the entropy.

Ans 3 possible outcomes with equal prob = $\frac{1}{3}$

$$\begin{aligned} (E) &= - \left(P_A \log_2 P_A + P_B \log_2 P_B + P_C \log_2 P_C \right) \\ &= - \left(\frac{1}{3} \log_2 \frac{1}{3} + \frac{1}{3} \log_2 \frac{1}{3} + \frac{1}{3} \log_2 \frac{1}{3} \right) \\ &= - \frac{1}{3} \left(\log_2 \frac{1}{3} + \log_2 \frac{1}{3} + \log_2 \frac{1}{3} \right) \\ &= \underline{\underline{-(-1.59)}} = \underline{\underline{1.59}} \end{aligned}$$

* In information theory, the information content (gain) is maximized, and entropy is minimized.

prob 2

Data for DT is given below :

ID	Age	Income	Student	CR	Buy car
1	Y	high	NO	fair	NO
2	Y	H	NO	good	NO
3	M	H	NO	f	Y
4	old	m	NO	f	Y
5	O	l	Y	f	Y
6	O	l	Y	g	N
7	m	l	Y	g	Y
8	Y	m	NO	f	N
9	Y	l	Y	f	Y
10	O	m	Y	f	Y
11	Y	m	Y	g	Y
12	M	m	NO	g	Y
13	M	h	Y	f	Y
14	O	m	NO	g	N

3 attribute seln. measures are

- (1) Information Gain
- (2) Gain Ratio (split info)
- (3) Gini Index

1) Information Gain :

$$Info(D) = - \sum_{i=1}^m P_i \log_2 P_i \quad \text{(also known as Entropy)} \quad - (1)$$

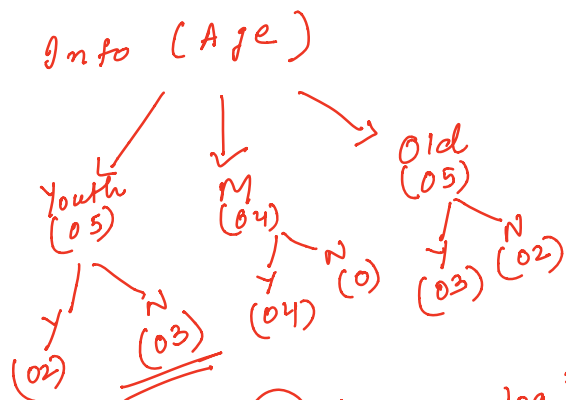
$$Info_A(D) = \sum_{j=1}^V \frac{|D_j|}{D} \times Info(D_j) \quad - (2)$$

$$Gain(A) = Info(D) - Info_A(D) \quad - (3)$$

Soln: 14 data set \rightarrow 9 (pos- buy car-yes)
 \rightarrow 5 (neg- buy car- No)

$$\text{Info (D)} = -\frac{9}{14} \log_2 \frac{9}{14} - \frac{5}{14} \log_2 \frac{5}{14}$$

$$= 0.94 \text{ bits}$$



$$\text{Info (Age)} = \frac{5}{14} \left(-\frac{2}{5} \log_2 \frac{2}{5} - \frac{3}{5} \log_2 \frac{3}{5} \right) +$$

$$\frac{4}{14} \left(-\frac{4}{4} \log_2 \frac{4}{4} \right) + \frac{5}{14} \left(-\frac{3}{5} \log_2 \frac{3}{5} - \frac{2}{5} \log_2 \frac{2}{5} \right)$$

$$= 0.694$$

$$\text{Gain (Age)} = \text{Info (D)} - \text{Info (Age)}$$

$$= 0.94 - 0.694$$

$$= \boxed{0.246} \quad \checkmark \quad (\text{maximum})$$

$$\text{Gain (Income)} = 0.029 \text{ bits}$$

$$\text{Gain (Student)} = 0.151 \text{ bits}$$

$$\text{Gain (CR)} = 0.048 \text{ bits}$$

\rightarrow Check for maximum information gain or min entropy value for every attribute, so as to select the splitting attribute