

ATTRIBUTE	POSSIBLE VALUES
age	old, midlife, new
competition	no, yes
type	software, hardware

AGE	COMPETITION	TYPE	PROFIT
old	yes	swr	down
old	no	swr	down
old	no	hwr	down
mid	yes	swr	down
mid	yes	hwr	down
mid	no	hwr	up
mid	no	swr	up
new	yes	swr	up
new	no	hwr	up
new	no	swr	up

$$\text{Gain}(X,T) = \text{Information}(X) - \text{Information}(X, T)$$

$$\text{Information}(X) = p(1) \cdot \log_2(p(1)) + \dots + p(n) \cdot \log_2(p(n))$$

where p_n is the fraction of the data set for which the value of the decision variable is the n th value. For instance $p(\text{down}) = 5/10$

$$\text{Information}(X,T) = \sum \text{over values of } T \left(\frac{|T_i|}{|T|} \cdot \text{Information}(T) \right)$$

So for Age: $\text{sum}(\text{old}, \text{mid}, \text{new}) \left(\frac{3}{10} \cdot \text{info}(\text{the 3 old}) + \frac{4}{10} \cdot \text{info}(\text{the 4 mid}) \dots \right)$

Algorithm

- Compute Information among current set of examples
- Compute Information Gain from splitting on all available features
- Add node to tree for splitting on the max gain.
- Split examples into subsets according to this new node in tree
- Return to step 1 with each subset