

Solution to Decision Trees Prac Sheet

Last modified on Wednesday, 24 April 2019 by f.maire@qut.edu.au

Exercise 1

- Consider the data table below

D =

X_1	X_2	X_3	X_4	C
F	F	F	F	P
F	F	T	T	P
F	T	F	T	P
T	T	T	F	P
T	F	F	F	N
T	T	T	T	N
T	T	T	F	N

$$\mathbf{X} = \{X_1, X_2, X_3, X_4\}$$

- What is the entropy of D ?

Write P_p for the probability of $C=P$, and P_n for the probability of $C=N$.

We have $P_p = 4/7$ and $P_n = 3/7$

$$H(D) = -P_n \log_2(P_n) - P_p \log_2(P_p)$$

Numerically

$$P_p = 0.5714$$

$$P_n = 0.4286$$

The entropy $H(D)$ is 0.9852

- What is the information gain of X_1 ?

$$H(C | X_1=N) = H([0-,3+]) = 0$$

$$H(C | X_1=P) = H([3-,1+]) = 0.8113$$

$$H(C | X_1) = 3/7 H([0-,3+]) + 4/7 H([3-,1+]) = 0.4636$$

$$\text{Information gain of } X_1 : H(D) - H(C | X_1) = 0.5216$$

- What is the information gain of X_2 ?

$$\text{Similarly, Gain}(X_2) = 0.98 - 0.97 = 0.01$$

- Build a DT to a depth of 3.

First split

D =

X ₁	X ₂	X ₃	X ₄	C
F	F	F	F	P
F	F	T	T	P
F	T	F	T	P
T	T	T	F	P
T	F	F	F	N
T	T	T	T	N
T	T	T	F	N

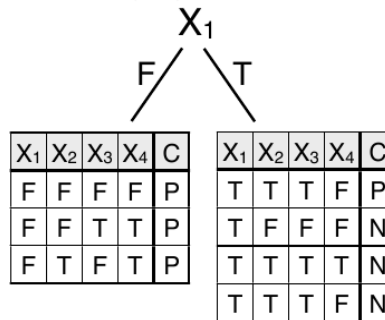
$$\mathbf{X} = \{X_1, X_2, X_3, X_4\}$$

$$\text{Gain}(X_1) = 0.52$$

$$\text{Gain}(X_2) = 0.01$$

$$\text{Gain}(X_3) = 0.01$$

$$\text{Gain}(X_4) = 0.01$$



Left subtree,

D =

X ₁	X ₂	X ₃	X ₄	C
F	F	F	F	P
F	F	T	T	P
F	T	F	T	P

$$\mathbf{X} = \{X_2, X_3, X_4\}$$

All instances have the same class.
Return class P.

Right subtree,

D =

X ₁	X ₂	X ₃	X ₄	C
T	T	T	F	P
T	F	F	F	N
T	T	T	T	N
T	T	T	F	N

X₂

F / \ T

0:1 1:2

X₃

F / \ T

0:1 1:2

X₄

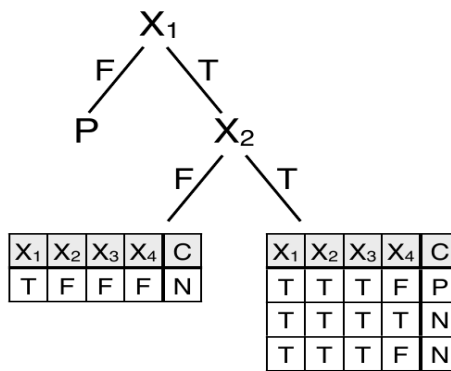
F / \ T

0:1 1:2

$$\mathbf{X} = \{X_2, X_3, X_4\}$$

All attributes have same information gain.
Break ties arbitrarily.
Choose X₂

So far we have,



For the left subtree of the tree rooted at X2, all instances have the same class, we will return class N.

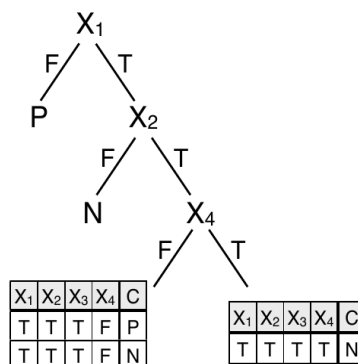
For the right subtree of the tree rooted at X2, we have

$$\mathbf{D} = \begin{array}{c|c|c|c|c|c} X_1 & X_2 & X_3 & X_4 & C \\ \hline T & T & T & F & P \\ \hline T & T & T & T & N \\ \hline T & T & T & F & N \end{array} \quad \begin{array}{c} X_3 \\ \hline F \quad T \\ \hline 0:0 \quad 1:2 \end{array} \quad \begin{array}{c} X_4 \\ \hline F \quad T \\ \hline 1:1 \quad 0:1 \end{array}$$

$\mathbf{X} = \{X_3, X_4\}$

X₃ has zero information gain
 X₄ has positive information gain
 Choose X₄

That is, we now have

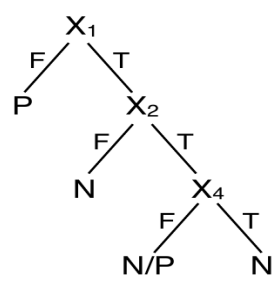


For the left subtree of X4, X3 has zero information gain. No suitable attribute for splitting.

Return most common class (break ties arbitrarily). Note: data is inconsistent!

For the right subtree of X4, All instances have the same class. Return N.

The final tree is



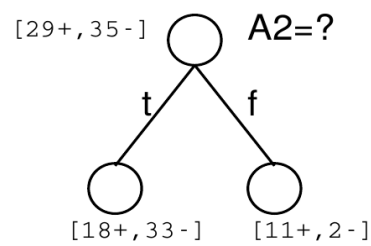
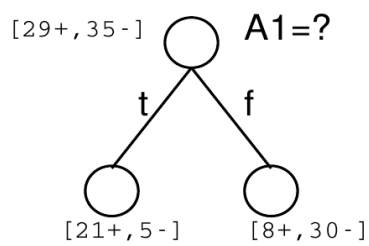
Exercise 2

- Recall that

$Gain(S, A) =$ expected reduction in entropy due to sorting on A

$$Gain(S, A) \equiv Entropy(S) - \sum_{v \in Values(A)} \frac{|S_v|}{|S|} Entropy(S_v)$$

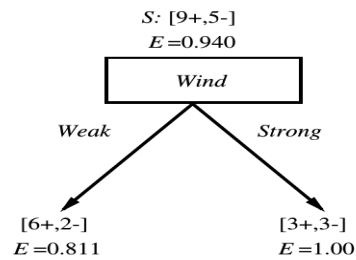
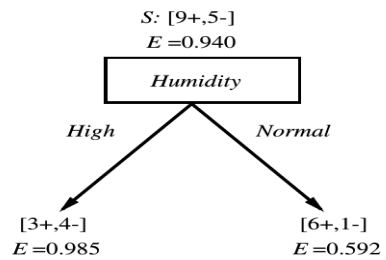
- Compute the information gain of the attribute $A1$ and $A2$



The information gain for A1 and A2 are respectively 0.2659 and 0.1214.

Exercise 3

- Which attribute is the best classifier?



We have

$$\begin{aligned} \text{Gain}(S, \text{Humidity}) &= .940 - (7/14).985 - (7/14).592 \\ &= .151 \end{aligned}$$

$$\begin{aligned} \text{Gain}(S, \text{Wind}) &= .940 - (8/14).811 - (6/14)1.0 \\ &= .048 \end{aligned}$$

Therefore Humidity is the most informative attribute with respect to class prediction.