

## CS 181 PROBLEM SET 1

ASHOK CUTKOSKY AND TONY FENG

**Problem 1.** (a) We calculate respective information content of  $A$  and  $B$ . First, we find the entropy of the variable.

$$H(X) = \left( \frac{4}{7} \log \frac{4}{7} + \frac{3}{7} \log \frac{3}{7} \right) \approx 0.985.$$

[This isn't actually necessary to determine which attribute has a higher information gain, unless they are equal and we are trying to determine if the gain is 0.]

- $A$  is true with probability  $\frac{4}{7}$ , with a 2, 2 split, and false with probability  $\frac{3}{7}$ , with a 2, 1 split. Therefore, the conditional entropy of  $A$  is

$$H(X | A) = - \left( \frac{4}{7} \left[ \frac{1}{2} \log \frac{1}{2} + \frac{1}{2} \log \frac{1}{2} \right] + \frac{3}{7} \left[ \frac{2}{3} \log \frac{2}{3} + \frac{1}{3} \log \frac{1}{3} \right] \right) \approx 0.965.$$

- $B$  is true with probability  $\frac{2}{7}$ , with a 1, 1 split, and false with probability  $\frac{5}{7}$ , with a 3, 2-split. Therefore,

$$H(X | B) = - \left( \frac{2}{7} \left[ \frac{1}{2} \log \frac{1}{2} + \frac{1}{2} \log \frac{1}{2} \right] + \frac{5}{7} \left[ \frac{3}{5} \log \frac{3}{5} + \frac{2}{5} \log \frac{2}{5} \right] \right) \approx 0.980.$$

Therefore,  $I(X, A) = H(X) - H(X | A) \geq H(X) - H(X | B) = I(X, B)$ . We conclude that  $A$  gives a higher information gain.

The differences and datasets are so small that it is difficult to give a convincing argument either way. The 2, 1 split when  $A$  is false perhaps suggests that it is a good indicator, but the sample size is just too small to be decisive. Furthermore, it appears that  $B$  does not split the data as well as  $A$ : its tree is more lopsided. A case for  $B$  could be made with the observation that the 2, 2 split when  $A$  is true suggests that it is a poor indicator; the 1, 1 split for  $B$  true is more noisy.

The example demonstrates that ID3's inductive bias that objects can be classified by sequential splitting on features one-by-one, unmindful of more "forward-looking" connections, can be flimsy in practice.

(b) Note that  $C$  and  $D$  have the exact same description of  $B$ , in terms of how many objects fell into each value of the attribute, and the distribution of classifications of each objects (in particular, each have one positive and one negative for some value of the attribute, and three positive and two negative for the other). Therefore, the calculation we did in (a) demonstrates that  $A$  has the highest information gain.

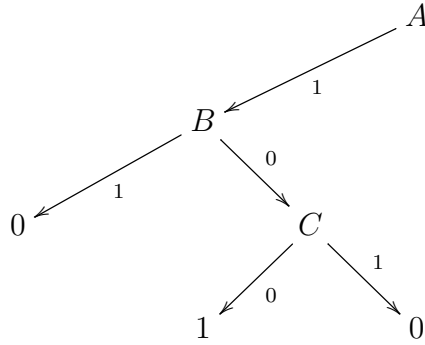
When  $A$  is true, the subdata obtained is

B	C	D	Label
0	1	0	0
0	0	0	1
0	0	1	1
1	0	1	0

Now  $B$  and  $C$  have identical classification results: one negative when 1, two positive and one negative when 0. On the other hand  $D$  evidently has zero mutual information: there is one negative and one positive in either case. Therefore,  $B$  and  $C$  clearly have equal mutual information with  $X$ , higher than  $D$ 's, and we split on  $B$  by alphabetical order. At this point, the subdata is

C	D	Label
1	0	0
0	0	1
0	1	1

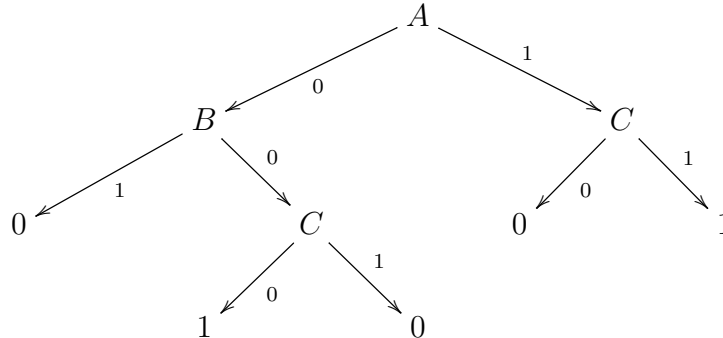
Now  $C$  perfectly classifies the data. To summarize, the left half of the tree we built is:



For the right half, the subdata is

B	C	D	Label
0	0	1	0
0	1	1	1
0	0	1	1

All examples have the same value for  $B$  and  $D$ , so clearly  $B$  and  $D$  offer no information whatsoever! Splitting on  $C$ , we find that there are two identical, inconsistent examples left, evenly split between 0 and 1. We arbitrarily choose to classify them as 1.



(c) Obviously, no tree can classify all examples since there are two inconsistent ones. However, some inspection shows that  $\text{not}(B \text{ xor } C)$  classifies all but the third example,  $0\ 0\ 0\ 1 \rightarrow 0$ .

