

COEN 240 MACHINE LEARNING

HOMEWROK FIVE

NAME: BOSEN YANG

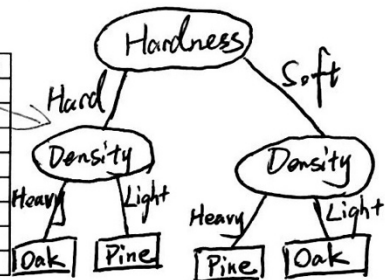
STUDENT ID: 1589880

Guideline: Please complete the following problems and generate a PDF file. Please refer to HomeworkFormat.pdf for the format of the submitted PDF file. *Therefore, Density will be chosen as node under condition Hardness = Soft. Till now, we've reached entropy = 0 on all leaf nodes. We are able to draw the complete decision tree as follows.*

Problem 1

You are a robot in a lumber yard, and must learn to discriminate Oak wood from Pine wood. You choose to learn a Decision Tree classifier. You are given the following examples:

| Example | Density | Grain | Hardness | Class |
|------------|---------|-------|----------|-------|
| Example #1 | Heavy | Small | Hard | Oak |
| Example #2 | Heavy | Large | Hard | Oak |
| Example #3 | Heavy | Small | Hard | Oak |
| Example #4 | Light | Large | Soft | Oak |
| Example #5 | Light | Large | Hard | Pine |
| Example #6 | Heavy | Small | Soft | Pine |
| Example #7 | Heavy | Large | Soft | Pine |
| Example #8 | Heavy | Small | Soft | Pine |



1.1 Which attribute will be chosen as the root of the tree (show derivations)?

1.2 Derive the complete decision tree by recursively applying the smallest entropy criterion to select root nodes of sub-trees (show derivations). Then draw the complete decision tree.

$$1.1 \quad H(C|Density) = \frac{6}{8} \cdot H(C|Density=Heavy) + \frac{2}{8} \cdot H(C|Density=Light) = \frac{6}{8} \cdot 1 + \frac{2}{8} \cdot 1 = 1 \text{ bit}$$

$$H(C|Grain) = \frac{4}{8} \cdot H(C|Grain=Small) + \frac{4}{8} \cdot H(C|Grain=Large) = \frac{4}{8} \cdot 1 + \frac{4}{8} \cdot 1 = 1 \text{ bit}$$

$$H(C|Hardness) = \frac{4}{8} \cdot H(C|Hardness=Hard) + \frac{4}{8} \cdot H(C|Hardness=Soft) = \frac{4}{8} \cdot 1 + \frac{4}{8} \cdot 1 = 1 \text{ bit}$$

Hardness will be chosen as the root because it has the smallest entropy (uncertainty).

$$1.2 \quad H(C|Density, Hardness=Hard) = \frac{3}{4} \cdot H(1,0) + \frac{1}{4} \cdot H(0,1) = 0 \text{ bit}$$

$$H(C|Grain, Hardness=Hard) = \frac{2}{4} \cdot H(\frac{3}{2}, \frac{1}{2}) + \frac{2}{4} \cdot H(0,1) > H(C|Density, Hardness=Hard)$$

$$\text{We choose Density as node under condition Hardness=Hard.}$$

$$H(C|Density, Hardness=Soft) = \frac{1}{4} \cdot H(1,0) + \frac{3}{4} \cdot H(0,1) = 0 \text{ bit}$$

$$H(C|Grain, Hardness=Soft) = \frac{2}{4} \cdot H(\frac{1}{2}, \frac{1}{2}) + \frac{2}{4} \cdot H(0,1) > H(C|Density, Hardness=Soft)$$

Problem 2

NASA wants to discriminate Martians (M) from Humans (H) based on these features (attributes): Green $\in \{N, Y\}$, Legs $\in \{2, 3\}$, Height $\in \{S, T\}$, Smelly $\in \{N, Y\}$. Your available training data is as follows (N=No, Y=Yes, S=Short, T=Tall):

| Example Number | Height | Green | Legs | Smelly | Target: Species |
|----------------|--------|-------|------|--------|-----------------|
| 1 | S | Y | 3 | Y | M |
| 2 | T | Y | 3 | N | M |
| 3 | S | Y | 3 | N | M |
| 4 | T | Y | 3 | N | M |
| 5 | T | N | 2 | Y | M |
| ?? 6 | T | Y | 2 | Y | H |
| 7 | S | N | 2 | N | H |
| 8 | T | N | 3 | N | H |
| 9 | S | N | 3 | N | H |
| 10 | T | N | 3 | N | H |

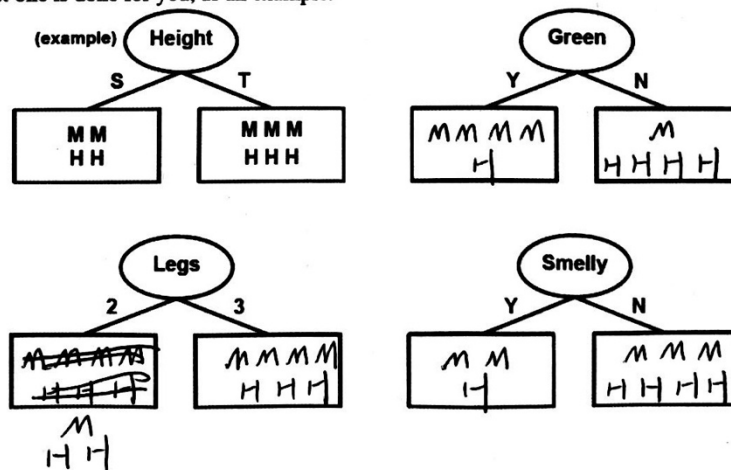
Please note:
A human might be green or have three legs for many possible reasons, e.g., if they were an actor playing a Martian as a role in a film or play. Anyway, it's a made-up problem

Great Point!

(a) What is the entropy of the target species before testing any attribute?

$$H\left(\frac{5}{10}, \frac{5}{10}\right) = H\left(\frac{1}{2}, \frac{1}{2}\right) = 1 \text{ bit}$$

(b) For each possible choice of root attribute below, show the resulting species distribution. Give your answer as M over H. The first one is done for you, as an example.



(c) What is the conditional entropy of species under attribute Height?

$$H(C|Height) = \frac{4}{10} \cdot H(\frac{1}{2}, \frac{1}{2}) + \frac{6}{10} \cdot H(\frac{1}{2}, \frac{1}{2}) = H(\frac{1}{2}, \frac{1}{2}) = 1 \text{ bit}$$

(d) What is the conditional entropy of species under attribute Green?

$$H(C|Green) = \frac{5}{10} \cdot H(\frac{4}{5}, \frac{1}{5}) + \frac{5}{10} \cdot H(\frac{1}{5}, \frac{4}{5}) = H(\frac{1}{5}, \frac{4}{5}) = 0.722 \text{ bit}$$

(e) What is the conditional entropy of species under attribute Legs?

$$H(C|Legs) = \frac{3}{10} \cdot H(\frac{1}{3}, \frac{2}{3}) + \frac{7}{10} \cdot H(\frac{4}{7}, \frac{3}{7}) = \frac{3}{10} \cdot 0.9183 + \frac{7}{10} \cdot 0.9852 = 0.965 \text{ bit}$$

(f) What is the conditional entropy of species under attribute Smelly?

$$H(C|Smelly) = \frac{3}{10} \cdot H(\frac{2}{3}, \frac{1}{3}) + \frac{7}{10} \cdot H(\frac{3}{7}, \frac{4}{7}) = H(C|Legs) = 0.965 \text{ bit}$$

(g) Which attribute would you select as the root attribute (i.e., the attribute to test first)? Why?

Green will be selected because it has the smallest entropy/uncertainty.

(h) Derive the complete decision tree by recursively applying the smallest entropy criterion to select root nodes of sub-trees (show derivations). Then draw the complete decision tree.

From previous questions we know Green will be chosen as root.

$$\left. \begin{aligned} H(C|Height, Green=Y) &= \frac{2}{5} \cdot H(1, 0) + \frac{3}{5} \cdot H(\frac{2}{3}, \frac{1}{3}) \\ H(C|Legs, Green=Y) &= \frac{4}{5} \cdot H(1, 0) + \frac{1}{5} \cdot H(0, 1) \end{aligned} \right\} \text{Legs will be selected as node under condition Green=Y.}$$

$$\left. \begin{aligned} H(C|Height, Green=N) &= \frac{2}{5} \cdot H(0, 1) + \frac{3}{5} \cdot H(\frac{1}{3}, \frac{2}{3}) \\ H(C|Legs, Green=N) &= \frac{3}{5} \cdot H(0, 1) + \frac{2}{5} \cdot H(\frac{1}{2}, \frac{1}{2}) \\ H(C|Smelly, Green=N) &= \frac{1}{5} \cdot H(1, 0) + \frac{4}{5} \cdot H(0, 1) \end{aligned} \right\} \text{Smelly will be selected as node under condition Green=N.}$$

Now on all leaf nodes we have entropy = 0, which means we are able to draw the complete decision tree.

