

AI PS5

Saturday, March 27, 2021 4:07 PM

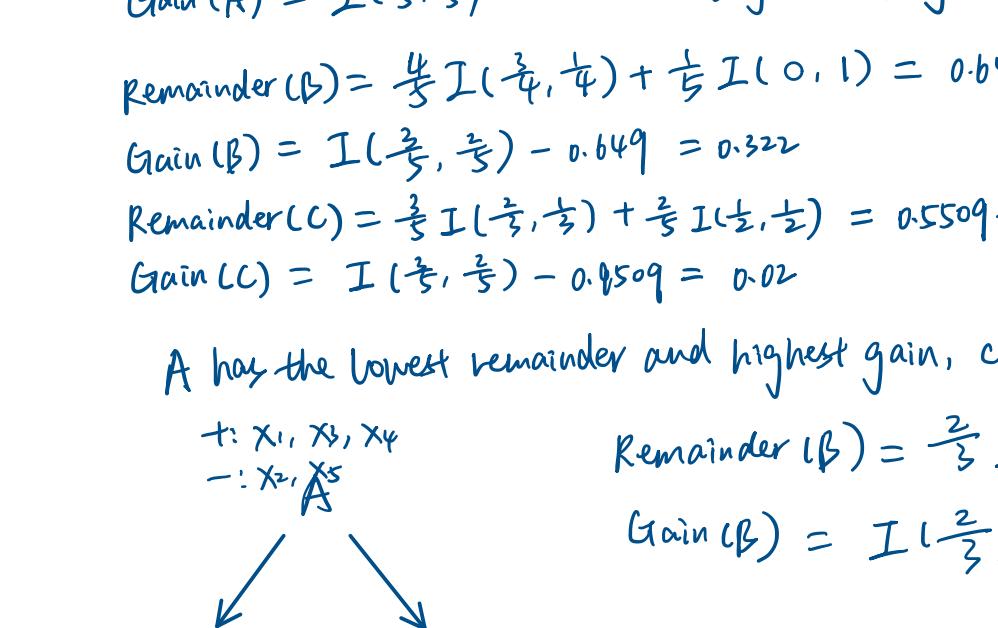
Problem #1 : Decision Trees (6 Points)

Given the dataset below that consists of five examples (X_1 through X_5) for three Boolean attributes (A, B, and C) and a Boolean outcome, you are to assemble (by hand) an optimal decision tree for this data.

	A	B	C	Outcome
X_1	T	T	F	T
X_2	F	T	F	F
X_3	F	T	T	T
X_4	T	T	T	T
X_5	F	F	T	F

Your solution should show all calculations, starting with a comparison between using each of the three attributes (A, B, C) as the root of the tree. You should show the Remainder values for each comparison, and show the final decision tree that represents the optimal decision tree.

Please submit your solution to Gradescope. Your solution will be scored based on the accuracy of your calculations and the clarity of your explanations.



Since an attribute A with d distinct values divides the training set E into subsets E_1, \dots, E_d . Each E_k has p_k positive examples and n_k negative examples. Then

$$\text{Remainder}(A) = \sum_{i=1}^d \frac{p_i + n_i}{p + n} I\left(\frac{p_i}{p + n}, \frac{n_i}{p + n}\right), \text{ where } I(p(v_1), \dots, p(v_n)) = \sum_{i=1}^n p(v_i) \log_2 p(v_i).$$

$$\text{Gain}(A) = I\left(\frac{p}{p+n}, \frac{n}{p+n}\right) - \text{Remainder}(A)$$

$$\text{Remainder}(A) = \frac{2}{5} I\left(\frac{2}{5}, \frac{3}{5}\right) + \frac{3}{5} I\left(\frac{1}{3}, \frac{2}{3}\right) = \frac{2}{5} \cdot \left(-\frac{1}{2} \log_2 \frac{1}{2} - \frac{3}{5} \log_2 \frac{3}{5}\right) = 0.6 \cdot 0.918 = 0.55$$

$$\text{Gain}(A) = I\left(\frac{2}{5}, \frac{3}{5}\right) - 0.55 = -\frac{2}{5} \log_2 \left(\frac{2}{5}\right) - \frac{3}{5} \log_2 \left(\frac{3}{5}\right) - 0.55 = 0.42$$

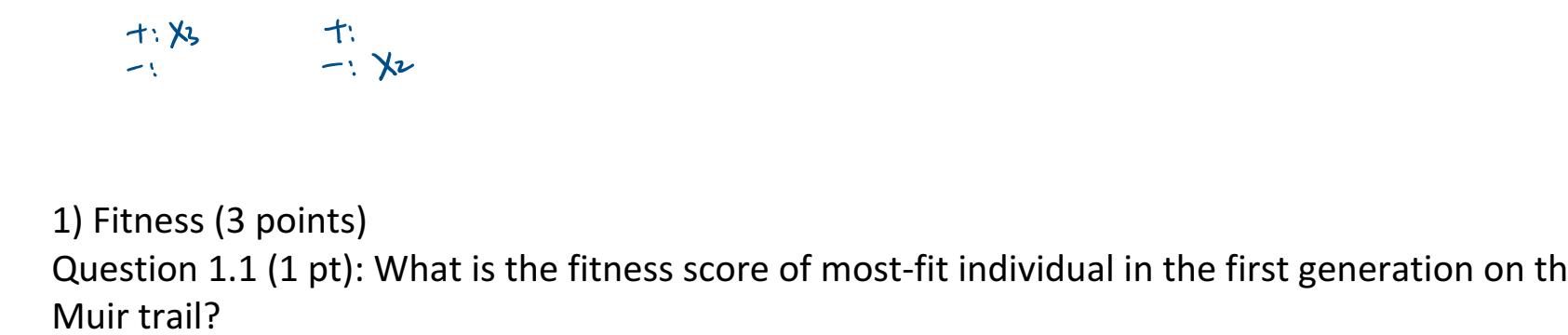
$$\text{Remainder}(B) = \frac{4}{5} I\left(\frac{1}{4}, \frac{3}{4}\right) + \frac{1}{5} I(0, 1) = 0.649$$

$$\text{Gain}(B) = I\left(\frac{2}{3}, \frac{1}{3}\right) - 0.649 = 0.322$$

$$\text{Remainder}(C) = \frac{2}{5} I\left(\frac{1}{2}, \frac{1}{2}\right) + \frac{3}{5} I\left(\frac{1}{2}, \frac{1}{2}\right) = 0.5509 + 0.4 = 0.9509$$

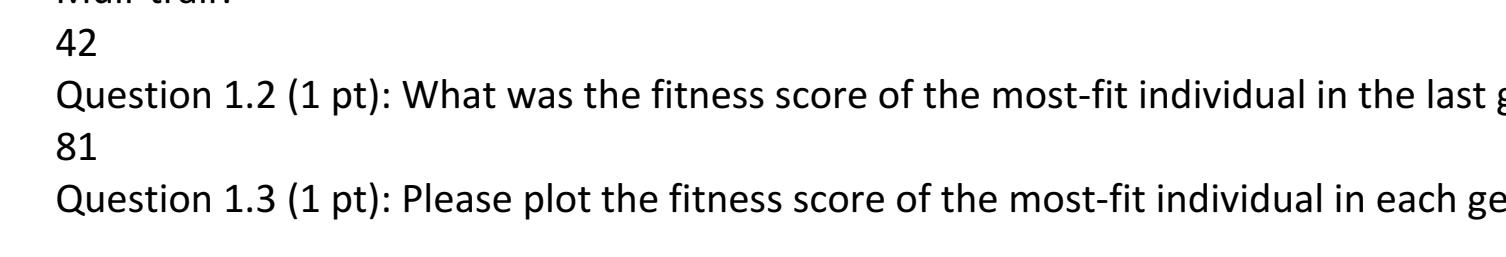
$$\text{Gain}(C) = I\left(\frac{1}{2}, \frac{1}{2}\right) - 0.9509 = 0.02$$

A has the lowest remainder and highest gain, choose A as root node.



$$\text{Remainder}(B) = \frac{2}{5} I\left(\frac{1}{2}, \frac{1}{2}\right) + \frac{1}{5} I(0, 1) = \frac{2}{5} \times 1 = 0.667$$

$$\text{Gain}(B) = I\left(\frac{2}{3}, \frac{1}{3}\right) - 0.667 = 0.251$$



$$\text{Remainder}(C) = \frac{2}{5} I\left(\frac{1}{2}, \frac{1}{2}\right) + \frac{1}{5} I(0, 1) = 0.667.$$

$$\text{Gain}(C) = I\left(\frac{1}{2}, \frac{1}{2}\right) - 0.667 = 0.251$$

B and C are equivalent!

Let's choose B. The optimal tree solution is

