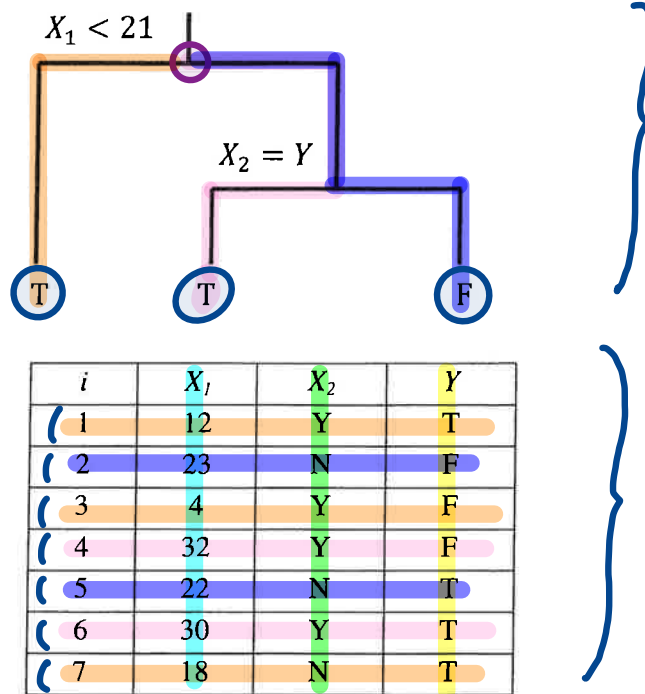


40.

You are given the following classification decision tree and data set:



Determine the relationship between the classification error rate, the Gini index, and the cross-entropy, summed across all nodes.

- A. cross-entropy > Gini index > classification error rate
- B. cross-entropy > Gini index = classification error rate
- C. classification error rate > Gini index > cross-entropy
- D. Gini index > cross-entropy > classification error rate
- E. The answer is not given by (A), (B), (C), or (D).

Caveat: They explicitly say: "summed across all nodes" which is different from computing a weighted average!

→:

For $X_1 < 21$, we have observations $i=1, 3, 7$ in that region. They have $Y_1 = \underline{T}$, $Y_3 = \underline{F}$, $Y_7 = \underline{T}$

⇒ From the tree, we know that the classification @ that node is **T**

⇒ The classification error is $\left(\frac{1}{3}\right)$

For $X_1 \geq 21$ and $X_2 = Y$, observations $i=4, 6$

are in that terminal node w/ $Y_4 = F$, $Y_6 = T$

⇒ The classification error is $\left(\frac{1}{2}\right)$

For $X_1 \geq 21$ and $X_2 = N$, observations $i=2, 5$

are in that terminal node w/ $Y_2 = F$, $Y_5 = T$

⇒ The classification error is $\left(\frac{1}{2}\right)$

The overall classification error: $\frac{1}{3} + \frac{1}{2} + \frac{1}{2} = \left(\frac{4}{3}\right) = \frac{12}{9}$

At the 1st terminal node: the Gini index = $\frac{\hat{p}}{3} \cdot (1 - \frac{\hat{p}}{3}) + \frac{2}{3} (1 - \frac{2}{3})$

At the 2nd terminal node: $2 \cdot \frac{1}{2} \cdot \frac{1}{2} = \frac{1}{2}$

At the 3rd — || — : the same

The total Gini index: $\frac{4}{9} + \frac{1}{2} + \frac{1}{2} = \left(\frac{13}{9}\right)$

The cross entropy @ 1st node: $-\frac{1}{3} \ln\left(\frac{1}{3}\right) - \frac{2}{3} \ln\left(\frac{2}{3}\right)$

The cross entropy @ 2nd and 3rd nodes:

$$-\frac{1}{2} \ln\left(\frac{1}{2}\right) - \frac{1}{2} \ln\left(\frac{1}{2}\right)$$

The total Cross Entropy: 2.022809

