

1. At root node, what is the Information Gain if we split by 'Age', 'Gender', and 'Car Ownership' respectively? (10 points)

HW 1. 100

① Split "Age": young, mid-age, old
old: Y, N, Y

$$H(S) = -p(+) \log_2 p(+) - p(-) \log_2 p(-)$$

$$\text{Entropy}[S, \text{Age}] = -\frac{3}{8} \log_2 \frac{3}{8} - \frac{5}{8} \log_2 \frac{5}{8} = 0.9544$$

~~Gain (Age) = 0.466~~

old: Entropy $[2^+, 1^-] = -\frac{2}{3} \log_2 \frac{2}{3} - \frac{1}{3} \log_2 \frac{1}{3} = 0.9183$

mid: Entropy $[1^+, 1^-] = 1$

young: Entropy $[0^+, 3^-] = 0$

$$\text{Gain}(S, \text{Age}) = 0.9544 - \left(\frac{3}{8} \cdot 0.9183\right) - \left(\frac{2}{8} \cdot 1\right) - \left(\frac{3}{8} \cdot 0\right)$$

$$= 0.56$$

② split "Gender": M, F

M: Entropy $[3^+, 5^-] = 0.9544$

$$\text{Gain}(S, \text{Gender}) = 0.9544 - \frac{8}{8} \cdot 0.9544 = 0$$

③ split "Car-Ownership": Y, N

Y: Entropy $[3^+, 2^-] = -\frac{3}{5} \log_2 \frac{3}{5} - \frac{2}{5} \log_2 \frac{2}{5} = 0.971$

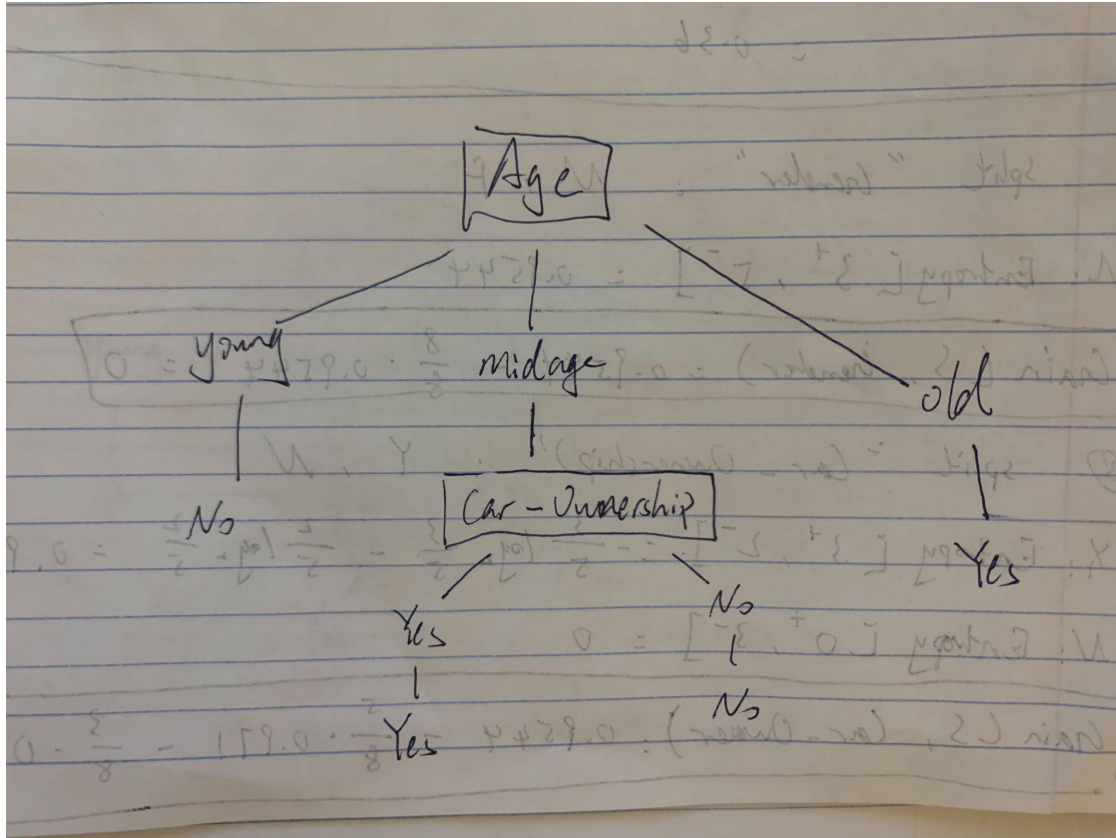
N: Entropy $[0^+, 3^-] = 0$

$$\text{Gain}(S, \text{Car-Owner}) = 0.9544 - \frac{5}{8} \cdot 0.971 - \frac{3}{8} \cdot 0 = 0.3475$$

2. Which column is used to make the first split? (10 points)

Age will be the column that is used to make the first split, because we need to choose the highest information gain to be our first split

3. Once trained, What will be the output on a new sample with Age = 50, Car Ownership = No, Gender = M? Can you train the same tree without the Gender column, why or why not? (10 points)



The output should be Yes, because when we put [Age = 50, Car Ownership = No, Gender = M], we will first go to the old branch, since we removed all the other attributes in the old branch. Then we can directly go to output Yes.

I think we can train the same tree without a Gender column, because all of the genders are Male, which makes no impact and no influence to the tree. Also, based on the information gain, we know that the information gain for Gender is 0. Since it is too low, we can just drop that.

The calculation of the decision tree: (next page)

Decision tree calculation

Car owner: old

~~Entropy~~ ~~old~~

$$\text{Gain (old, Car-owner)} = 0.91 - \frac{3}{3} \cdot 0.9185 = 0$$

$$\text{Gain (old, Gender)} = 0.91 - \frac{3}{2} \cdot 0.9185 = 0$$

$$\text{Entropy (old)} = -\frac{2}{3} \log_2 \frac{2}{3} - \frac{1}{3} \log_2 \frac{1}{3} = 0.9185$$

$$\text{Entropy (old, car owner)} = 0.9185$$

$$\text{Gain (mid, Car-owner)} = 1 - \frac{1}{2} \cdot 0 - \frac{1}{2} \cdot 0 = 1$$

$$\text{Yes: } E[1^+, 0^-] = 0$$

$$\text{No: } E[0, 1^-] = 0$$

$$\text{Gain (mid, Gender)} = 1 - \frac{2}{2} \cdot 1 = 0$$

$$\text{Male: } E[1^+, 1^-] = 1$$

$$\text{Gain (young, Car)} = 0 - \frac{1}{3} \cdot 0 - \frac{2}{3} \cdot 0 = 0$$

$$\text{Yes: } E[0^+, 1] = 0$$

$$\text{No: } E[0^+, 2^-] = 0$$

$$\text{Gain (young, Gender)} = 0 - \frac{3}{3} \cdot 0 = 0$$

$$\text{No: } E[0^+, 5^-] = 0$$