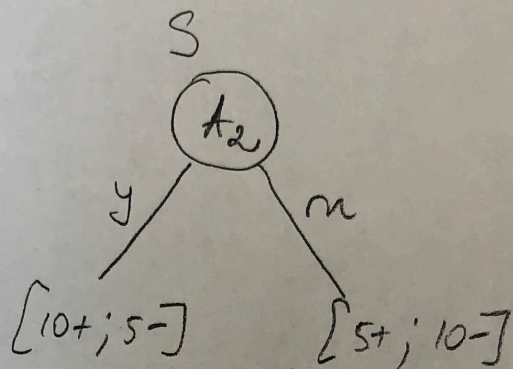
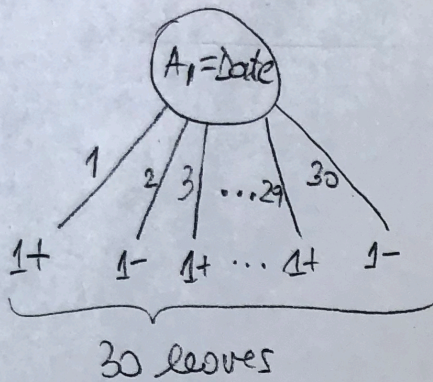


Example: Gain Ratio (4) vs. Info. gain (2)

$$S = 30 \text{ ex.} = [15+, 15-]$$



perfectly predicts training data but poorly the test data

$$1) E(S) = -\frac{15}{30} \log_2 \frac{15}{30} - \frac{15}{30} \log_2 \frac{15}{30} = 1$$

$$2) \text{Gain}(S, \text{Date}) = 1 - \underbrace{\frac{1}{30}(0) - \frac{1}{30} \cdot 0 - \dots - \frac{1}{30} \cdot 0}_{30 \text{ times}} = 1$$

$$\text{Gain}(S, A_2) = 1 - \frac{15}{30} \left( -\frac{10}{15} \log_2 \frac{10}{15} - \frac{5}{15} \log_2 \frac{5}{15} \right) - \frac{15}{30} \left( -\frac{5}{15} \log_2 \frac{5}{15} - \frac{10}{15} \log_2 \frac{10}{15} \right) = 0.8617$$

Because  $1 > 0.8617 \xrightarrow{\text{Info Gain}} A_1 = \text{Date}$  will be selected  $\rightarrow$  NOT GOOD!

$$3) \text{Split Info}(S, \text{Date}) = \underbrace{-\frac{1}{30} \log_2 \frac{1}{30} - \dots - \frac{1}{30} \log_2 \frac{1}{30}}_{30 \text{ times}} = -\log_2 \frac{1}{30} = 4.9069$$

$$\text{Split Info}(S, A_2) = -\frac{15}{30} \log_2 \frac{15}{30} - \frac{15}{30} \log_2 \frac{15}{30} = -\log_2 \frac{1}{2} = 1$$

$$4) \text{Gain Ratio}(S; \text{Date}) = \frac{1}{4.9069} = 0.2038 < \text{Gain Ratio}(S; A_2) = \frac{0.8617}{1} = 0.8617$$