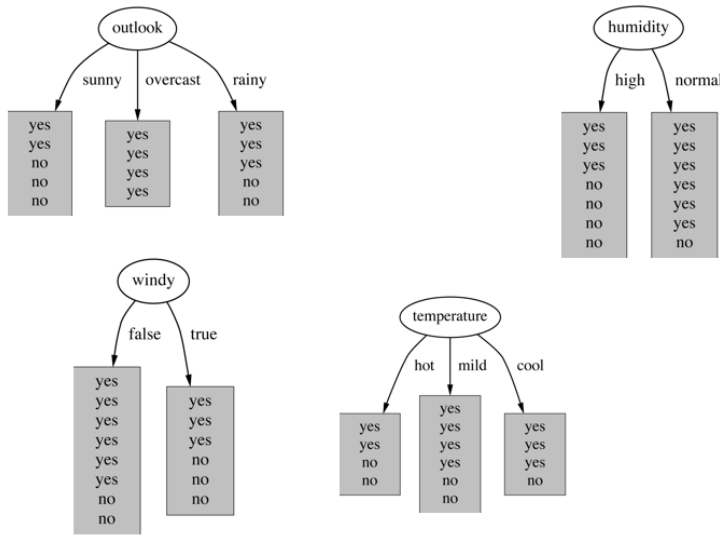


Data Mining, CSCI 347, Fall 2019
Decision Trees, Sept. 18

A simple decision tree algorithm is to calculate the information gained when splitting on each attribute, split on the attribute that yields the greatest gain, and recursively repeat this process for each branch.

Consider the weather dataset and splitting on the outlook attribute.



1. Calculate the information of the top node.

$$\begin{aligned}
 \text{info}([9,5]) &= \text{entropy}(9/14, 5/14) \\
 &= -9/14 (\log_2 9/14) - 5/14 (\log_2 5/14) \\
 &= -9/14 (\log_{10} 9/14) / (\log_{10} 2) - 5/14 (\log_{10} 5/14) / (\log_2 2) \\
 &= 0.940 \text{ bits}
 \end{aligned}$$

2. Give the formula to calculate the information gained by splitting on Outlook.

$$\text{gain}(\text{Outlook}) = \text{info}([9,5]) - \text{info}([2,3], [4,0], [3,2])$$

3. Calculate the purity of each of the three branches.

Outlook = sunny

$$\begin{aligned} \text{info}([2,3]) &= \text{entropy}(2/5, 3/5) \\ &= -2/5 (\log_2 2/5) - 3/5 (\log_2 3/5) \\ &= -2/5 (\log_{10} 2/5) / (\log_{10} 2) - 3/5 (\log_{10} 3/5) / (\log_{10} 2) \\ &= 0.971 \text{ bits} \end{aligned}$$

Outlook = overcast

$$\begin{aligned} \text{info}([4,0]) &= \text{entropy}(4/4, 0/4) \\ &= \text{entropy}(1,0) \\ &= -1 (\log_2 1) - 0 (\log_2 0) \\ &\quad \text{Ordinarily } \log_2 0 \text{ is not defined, ok since times 0} \\ &= -1 (\log_{10} 1) / (\log_{10} 2) - 0 \\ &= -1 (0 / (\log_{10} 2)) - 0 \\ &= 0 \text{ bits} \end{aligned}$$

Outlook = rainy

$$\begin{aligned} \text{info}([3,2]) &= \text{entropy}(3/5, 2/5) \\ &= \text{entropy}(2/5, 3/5) \\ &= -2/5 (\log_2 2/5) - 3/5 (\log_2 3/5) \\ &= -2/5 (\log_{10} 2/5) / (\log_{10} 2) - 3/5 (\log_{10} 3/5) / (\log_{10} 2) \\ &= 0.971 \text{ bits} \end{aligned}$$

4. Calculate the weighted average information value for splitting on Outlook

$$\begin{aligned} \text{info}([3,2], [4,0], [3,2]) &= 5/14 * 0.971 + 4/14 * 0 + 5/14 * 0.971 \\ &= 0.693 \text{ bits} \end{aligned}$$

5. Information gain from splitting on Outlook.

$$\begin{aligned} \text{info}([9,5]) - \text{info}([3,2], [4,0], [3,2]) &= 0.940 - 0.693 \\ &= 0.247 \text{ bits} \end{aligned}$$

6. Calculate the intrinsic information value of the split

$$\begin{aligned} \text{info}([5,4,5]) &= \text{entropy}(5/14, 4/14, 5/14) \\ &= -5/14 (\log_2 5/14) - 4/14 (\log_2 4/14) - 5/14 (\log_2 5/14) \\ &= -5/14 (\log_{10} 5/14) / (\log_{10} 2) - 4/14 (\log_{10} 4/14) / (\log_{10} 2) \\ &\quad - 5/14 (\log_{10} 5/14) / (\log_{10} 2) \\ &= 1.577 \text{ bits} \end{aligned}$$

7. Calculate the ratio by dividing the attribute gain by the intrinsic information value

$$\begin{aligned} \text{Gain Ratio} &= \text{Gain from Attribute} / \text{Intrinsic Value of Split} \\ &= 0.247/1.577 \\ &= 0.157 \end{aligned}$$

8. Repeat the process the other attributes.

$$\begin{aligned} \text{Gain ratio} &= \text{information gain of attribute split} / \text{intrinsic information of attribute split} \\ &= (\text{information at root} - \text{average information for attribute split}) \\ &\quad / \text{intrinsic information of attribute split} \end{aligned}$$

$$\begin{aligned} &\text{Gain ratio for splitting on Outlook} \\ &= (\text{information at root} - \text{average information for attribute split}) \\ &\quad / \text{intrinsic information of attribute split} \\ &= (\text{info}([9,5]) - \text{info}([2,3],[4,0],[3,2])) / \text{info}([5,4,5]) \\ &= (\text{info}([9,5]) - (5/14*\text{info}([2,3]) + 4/14*\text{info}([4,0]) + 5/14*\text{info}([3,2]))) \\ &\quad / \text{info}([5,4,5]) \\ &= (\text{entropy}(9/14,5/14) - \\ &\quad (5/14*\text{entropy}(2/5,3/5) + 4/14*\text{entropy}(4/4,0/4) + \\ &\quad 5/14*\text{entropy}(3/5,2/5)) \\ &\quad / \text{entropy}(5/14,4/14,5/14)) \\ &= (\text{entropy}(9/14,5/14) - \\ &\quad (5/14*\text{entropy}(2/5,3/5) + 4/14*\text{entropy}(4/4,0/4) + \\ &\quad 5/14*\text{entropy}(3/5,2/5)) \\ &\quad / \text{entropy}(5/14,4/14,5/14)) \\ &= (-9/14*\log_2(9/14) - 5/14*\log_2(5/14)) - \\ &\quad (5/14*[-2/5*\log_2(2/5) - 3/5*\log_2(3/5)] \\ &\quad + 4/14*[-4/4*\log_2(4/4) - 0/4*\log_2(0/4)] \\ &\quad + 5/14*[-3/5*\log_2(3/5) - 2/5*\log_2(2/5)]) \\ &\quad / (-5/14*\log_2(5/14) - 4/14*\log_2(4/14) - 5/14*\log_2(5/14)) \\ &= (0.940 \text{ bits} - \\ &\quad (5/14*0.971 \text{ bits} + 4/14* 0 \text{ bits} + 5/14 * 0.971 \text{ bits})) \\ &\quad / 1.577 \text{ bits} \\ &= (0.940 \text{ bits} - \\ &\quad (0.347 \text{ bits} + 0 \text{ bits} + 0.347 \text{ bits})) \\ &\quad / 1.577 \text{ bits} \\ &= (0.940 \text{ bits} - 0.694 \text{ bits}) / 1.577 \text{ bits} \\ &= 0.246 \text{ bits} / 1.577 \text{ bits} \\ &= 0.156 \end{aligned}$$

Gain ratio for splitting on Temperature

$$\begin{aligned}
 &= (\text{information at root} - \text{average information for attribute split}) \\
 &\quad / \text{intrinsic information of attribute split} \\
 &= (\text{info}([9,5]) - \text{info}([2,2],[4,2],[3,1])) / \text{info}([4,6,4]) \\
 &= (\text{info}([9,5]) - (4/14 * \text{info}([2,2]) + 6/14 * \text{info}([4,2]) + 4/14 * \text{info}([3,1]))) \\
 &\quad / \text{info}([4,6,4]) \\
 &= (\text{entropy}(9/14,5/14) - \\
 &\quad (4/14 * \text{entropy}(2/4,2/4) + 6/14 * \text{entropy}(4/6,2/6) + \\
 &4/14 * \text{entropy}(3/4,1/4)) \\
 &\quad / \text{entropy}(4/14,6/14,4/14) \\
 &= (\text{entropy}(9/14,5/14) - \\
 &\quad (4/14 * \text{entropy}(2/4,2/4) + 6/14 * \text{entropy}(4/6,2/6) + \\
 &4/14 * \text{entropy}(3/4,1/4)) \\
 &\quad / \text{entropy}(4/14,6/14,4/14) \\
 &= (0.940 \text{ bits} - (4/14 * 1 + 6/14 * 0.438 + 4/14 * 0.791)) / 1.362 \text{ bits} \\
 &= (0.940 \text{ bits} - 0.911 \text{ bits}) / 1.362 \text{ bits} \\
 &= 0.029 \text{ bits} / 1.577 \text{ bits} \\
 &= 0.021
 \end{aligned}$$

Gain ratio for splitting on Humidity

$$\begin{aligned}
 &= (\text{information at root} - \text{average information for attribute split}) \\
 &\quad / \text{intrinsic information of attribute split} \\
 &= (\text{info}([9,5]) - \text{info}([3,4],[6,1])) / \text{info}([7,7]) \\
 &= (\text{info}([9,5]) - (7/14 * \text{info}([3,4]) + 7/14 * \text{info}([6,1]))) \\
 &\quad / \text{info}([7,7]) \\
 &= (\text{entropy}(9/14,5/14) - \\
 &\quad (7/14 * \text{entropy}(3/7,4/7) + 7/14 * \text{entropy}(6/7,1/7)) \\
 &\quad / \text{entropy}(7/14,7/14) \\
 &= (\text{entropy}(9/14,5/14) - \\
 &\quad (7/14 * \text{entropy}(3/7,4/7) + 7/14 * \text{entropy}(6/7,1/7)) \\
 &\quad / \text{entropy}(7/14,7/14) \\
 &= (0.940 \text{ bits} - 0.788 \text{ bits}) / 1 \text{ bit} \\
 &= 0.152 \text{ bits} / 1 \text{ bit} \\
 &= 0.152
 \end{aligned}$$

Gain ratio for splitting on Windy

$$\begin{aligned}
 &= (\text{information at root} - \text{average information for attribute split}) \\
 &\quad / \text{intrinsic information of attribute split} \\
 &= (\text{info}([9,5]) - \text{info}([6,2],[3,3])) / \text{info}([8,6]) \\
 &= (\text{info}([9,5]) - (8/14 * \text{info}([6,2]) + 6/14 * \text{info}([3,3]))) \\
 &\quad / \text{info}([8,6]) \\
 &= (\text{entropy}(9/14,5/14) - \\
 &\quad (8/14 * \text{entropy}(6/8,2/8) + 6/14 * \text{entropy}(3/6,3/6)) \\
 &\quad / \text{entropy}(8/14,6/14) \\
 &= (\text{entropy}(9/14,5/14) - \\
 &\quad (8/14 * \text{entropy}(6/8,2/8) + 6/14 * \text{entropy}(6/7,1/7))
 \end{aligned}$$

$$\begin{aligned} & / \text{entropy}(8/14,6/14) \\ & = (0.940 \text{ bits} - 0.892 \text{ bits}) / 0.985 \text{ bits} \\ & = 0.048 \text{ bits} / 0.985 \text{ bits} \\ & = 0.049 \end{aligned}$$

9. Based on these calculations, what attribute would appear as the root of the tree?

Outlook because it creates the greatest gain ratio.