

CIS 419 – Homework 1

Devin Stein
devstein@seas.upenn.edu

September 25, 2016

Problem 1

- a) At the root node for a decision tree in this domain, what are the information gains associated with the Outlook and Humidity attributes? (Use a threshold of 75 for humidity (i.e., assume a binary split: humidity ≤ 75 / humidity > 75). Be sure to show your computations.

$$\begin{aligned} H(\text{Root}) &= -\frac{9}{14} \log\left(\frac{9}{14}\right) - \frac{5}{14} \log\left(\frac{5}{14}\right) = .94 \\ H(\text{Humidity}) &= \frac{5}{14} \left(-\frac{4}{5} \log\left(\frac{4}{5}\right) - \frac{1}{5} \log\left(\frac{1}{5}\right) \right) + \frac{9}{14} \left(-\frac{5}{9} \log\left(\frac{5}{9}\right) - \frac{4}{9} \log\left(\frac{4}{9}\right) \right) = .895 \\ H(\text{Outlook}) &= \frac{5}{14} \left(-\frac{2}{5} \log\left(\frac{2}{5}\right) - \frac{3}{5} \log\left(\frac{3}{5}\right) \right) + \left(\frac{4}{14}(0) \right) + \frac{5}{14} \left(-\frac{3}{5} \log\left(\frac{3}{5}\right) - \frac{2}{5} \log\left(\frac{2}{5}\right) \right) = .694 \end{aligned}$$

$$\text{InformationGain}(\text{Humidity}) = .94 - .895 = .045$$

$$\text{InformationGain}(\text{Outlook}) = .94 - .694 = .247$$

- b) Again at the root node, what are the gain ratios associated with the Outlook and Humidity attributes (using the same threshold as in (a))? Be sure to show your computations.

$$\begin{aligned} \text{SplitInfo}(\text{Humidity}) &= -\frac{9}{14} \log\left(\frac{9}{14}\right) - \frac{5}{14} \log\left(\frac{5}{14}\right) = .94 \\ \text{SplitInfo}(\text{Outlook}) &= -\frac{5}{14} \log\left(\frac{5}{14}\right) - \frac{4}{14} \log\left(\frac{4}{14}\right) - \frac{5}{14} \log\left(\frac{5}{14}\right) = 1.58 \end{aligned}$$

$$\text{InformationRatio}(\text{Humidity}) = \frac{.045}{.94} = .048$$

$$\text{InformationRatio}(\text{Outlook}) = \frac{.247}{1.58} = .156$$

- c) Draw the complete (unpruned) decision tree, showing the class predictions at the leaves.

$\text{InformationRatio}(\text{Sunny}, \text{Humidity}) > \text{InformationRatio}(\text{Sunny}, \text{Wind}) > \text{InformationRatio}(\text{Sunny}, \text{Temp})$

$\text{InformationRatio}(\text{Rain}, \text{Wind}) > \text{InformationRatio}(\text{Rain}, \text{Humidity}) > \text{InformationRatio}(\text{Rain}, \text{Temp})$

Splitting on (Sunny, Humidity) gives the most information on a perfect Play/Don't Play split. Similarly, Splitting (Rain, Wind) gives is perfect. Lastly, only looking at (Overcast) tells you to

always play.







