

SMAI In class Assignment

Author: Vishal Reddy Mandadi

Roll Number: 2019101119

Question 1

Information gain for split 1

$$\begin{aligned}E(\text{parent}) &= - \sum (p(w_i) \log(p(w_i))) \\E(\text{parent}) &= -1 * 4 * \frac{10}{40} * \log_2(10/40) = -\log_2(1/4) = 2 \\E(\text{child1}) &= 2 * \frac{-10}{20} \log_2(10/20) = 1 \\E(\text{child2}) &= 2 * \frac{-10}{20} \log_2(10/20) = 1 \\ \text{information gain} &= i_{\text{parent}} - \frac{n_1}{n} i_{\text{child1}} - \frac{n_2}{n} i_{\text{child2}} \\ \Rightarrow \text{information gain} &= 2 - \left(\frac{20}{40} * 1\right) - \left(\frac{20}{40} * 1\right) = 2 - 1 = 1 \\ \therefore \text{information gain for split 1} &= 1\end{aligned}$$

Information gain for split 2

Assumption: Split 2 isn't considered as the split following split 1. It was mentioned by the professor that we need to consider split 2 as a complete 1st level root split (similarly split 3). This implies, split 2 included the split created by the straight line and the dotted line that continues up from the straight line

$$\begin{aligned}E(\text{parent}) &= - \sum (p(w_i) \log(p(w_i))) \\E(\text{parent}) &= -1 * 4 * \frac{10}{40} * \log_2(10/40) = -\log_2(1/4) = 2 \\E(\text{child1}) &= -1 * \left(\frac{10}{16} * \log_2\left(\frac{10}{16}\right) + \frac{6}{16} * \log_2\left(\frac{6}{16}\right)\right) \\&= 3 - \frac{5}{8} \log_2(5) - \frac{3}{8} \log_2(3) = 3 - 1.451 - 0.594 = 0.955 \\E(\text{child2}) &= -1 * \left(2 * \frac{10}{24} * \log_2\left(\frac{10}{24}\right) + \frac{4}{24} * \log_2\left(\frac{4}{24}\right)\right) \\E(\text{child2}) &= -1(0.833 * (-1.265) + (0.167 * (-2.585))) = 1.485 \\ \text{information gain} &= i_{\text{parent}} - \frac{n_1}{n} i_{\text{child1}} - \frac{n_2}{n} i_{\text{child2}} \\ \Rightarrow \text{information gain} &= 2 - \left(\frac{16}{40} * 0.955\right) - \left(\frac{24}{40} * 1.485\right) = 2 - 0.382 - 0.891 = 0.727 \\ \therefore \text{information gain for split 2} &= 0.727\end{aligned}$$

Information gain for split 3

$$\begin{aligned}
 E(\text{parent}) &= - \sum (p(w_i) \log(p(w_i))) \\
 E(\text{parent}) &= -1 * 4 * \frac{10}{40} * \log_2(10/40) = -\log_2(1/4) = 2 \\
 E(\text{child1}) &= -1 * \left(\frac{10}{14} * \log_2\left(\frac{10}{14}\right) + \frac{4}{14} * \log_2\left(\frac{4}{14}\right) \right) \\
 &= -1 * (0.714 * (-0.486) + 0.286 * (-1.806)) = 0.864 \\
 E(\text{child2}) &= -1 * \left(2 * \frac{10}{26} * \log_2\left(\frac{10}{26}\right) + \frac{6}{26} * \log_2\left(\frac{6}{26}\right) \right) \\
 E(\text{child2}) &= -1(2 * 0.385 * (-1.377) + (0.231 * (-2.114))) = 1.549 \\
 \text{information gain} &= i_{\text{parent}} - \frac{n_1}{n} i_{\text{child1}} - \frac{n_2}{n} i_{\text{child2}} \\
 \Rightarrow \text{information gain} &= 2 - \left(\frac{14}{40} * 0.864 \right) - \left(\frac{26}{40} * 1.549 \right) = 2 - 0.302 - 1.007 = 0.691 \\
 \therefore \text{information gain for split 2} &= 0.691
 \end{aligned}$$

Therefore, split 1 is most favorable. The order of favorable splits from most to least is

$$\text{split1} > \text{split2} > \text{split3}$$

Question 2

- a. The problem with the given data is that all possible first splits give 0 information gain, due to which we won't be able to find the most favorable first split. Thus the algorithm gets stuck at the first step itself and might take up a least favorable split as its first choice leading to erroneous classification.
- b. This problem can be solved in two ways:
 - a. Pre-processing the data using PCA (principle component analysis) - this way we will be able to get the data projected onto a new set of the axis where it can be easily split by decision trees
 - b. Using splits that are not parallel to the principle axis, i.e., using a linear combination of features instead of a single feature. The weights for the linear combination can be learned using gradient descent (similar to linear regression). The only downside with this method is that it trains very slowly