

ARIN - Statistical Machine Learning

ID3 Decision Tree Learning - Provisional Notes

shi at ebornet dot com

May 9, 2016

Disclaimer: This document is provided "as is". In no event shall the author be liable for damages resulting from or related to the use of this document.

1 Problem Scenarios

For a decision tree problem, there are two possible scenarios:

1. If there is a 'conclusion column' in the table, then that column's data will be our starting point – we do not need to choose a random variable with the highest initial entropy.
2. If there is no 'conclusion column', we must calculate the entropy of each random variable, and choose the variable with the highest entropy to establish the baseline.

For convenience, this note will assume a binary sample space (true or false), it is trivial to extend this summary of method to problems with non-binary sample space by calculating more branches for each node.

2 Compute Root Entropy

In Scenario 1, we only need to compute for the 'conclusion column' as stated by the question; while in Scenario 2, we need to compute for all possible random variables, and choose the random variable with the highest entropy for root entropy.

To compute the root entropy:

$$P(S_{initial}) = (\frac{true_count}{total_count}, \frac{false_count}{total_count})$$

$$\text{Entropy } H(S_{initial}) = -\frac{true_count}{total_count} \log_2 \frac{true_count}{total_count} - \frac{false_count}{total_count} \log_2 \frac{false_count}{total_count}$$

where *true_count* is the number of true results and *false_count* is the number of false results for the random variable, out of *total_count* which is the total number of samples.

3 Compute A Child Node

For each parent node (which may be the root node, depends on position in the tree), we need to compute entropy and information gain for each of its child (other random variables), until we have no more random variables.

For each child of a parent node:

$$Q_{child_true} = \frac{child_true_count}{parent_total}$$

$$Q_{child_false} = \frac{child_false_count}{parent_total}$$

where $child_true_count$ is the total number occurrences when the child random variable is true, regardless of the true or false value of its branch from the parent node. The same applies to $child_false_count$ when the child random variable is false.

$$parent_total = child_true_count + child_false_count$$

Then we have vectors $P(child_true)$ and $P(child_false)$ for both possibilities:

$$P(S_{child_true}) = (\frac{child_true_parent_true}{child_true_count}, \frac{child_true_parent_false}{child_true_count})$$

$$P(S_{child_false}) = (\frac{child_false_parent_true}{child_false_count}, \frac{child_false_parent_false}{child_false_count})$$

For convenience, we shall abbreviate the above terms as:

$$P(S_{child_true}) = (\frac{CTPT}{CTC}, \frac{CTPF}{CTC})$$

$$P(S_{child_false}) = (\frac{CFPT}{CFC}, \frac{CFPF}{CFC})$$

Now we can calculate the entropy for both possibilities of the child random variable:

$$\text{Entropy } H(S_{child_true}) = -\frac{CTPT}{CTC} \log_2 \frac{CTPT}{CTC} - \frac{CTPF}{CTC} \log_2 \frac{CTPF}{CTC}$$

$$\text{Entropy } H(S_{child_false}) = -\frac{CFPT}{CFC} \log_2 \frac{CFPT}{CFC} - \frac{CFPF}{CFC} \log_2 \frac{CFPF}{CFC}$$

With entropy for both possibilities, we are able to calculate the split:

$$H(Split) = \frac{child_true_count}{parent_total} \times H(S_{child_true}) + \frac{child_false_count}{parent_total} \times H(S_{child_false})$$

Finally, the information gain of the child node can be calculated from its parent's entropy and its split:

$$InformationGain = H(S_{parent}) - H(Split)$$

This process is repeated for all child nodes (remaining random variables) of the parent node. We choose the random variable with the highest *InformationGain* as the next best node in decision making. The next best node becomes the new parent node and this process continues, until exhaustion (as stated at the beginning of this section) or when both entropy values are 0 – the decision can then be made based on traversing the tree from the root node to that child node, and no more information is necessary.