

Information Gain Equation:

Importance = Original Entropy - New Entropy

$$\text{Original Entropy} = H\left(\frac{a}{a+b+c+\dots}, \frac{b}{a+b+c+\dots}, \frac{c}{a+b+c+\dots}, \dots\right)$$

Where a, b, c, \dots correspond to the number of remaining examples in the training data with some outcome and $a + b + c + \dots$ corresponds to the sum of all remaining examples in the training data at this point in the tree.

$$\text{New Entropy} = \sum_{k=1}^d \frac{a_k+b_k+c_k+\dots}{a+b+c+\dots} H\left(\frac{a_k}{a_k+b_k+c_k+\dots}, \frac{b_k}{a_k+b_k+c_k+\dots}, \frac{c_k}{a_k+b_k+c_k+\dots}, \dots\right)$$

Where d corresponds to the number of return values of the question and a_k, b_k, c_k, \dots correspond to the number of each outcome of the examples in the k return value of the question. $a + b + c + \dots$ is again the sum of all examples and $a_k + b_k + c_k + \dots$ is the sum of all examples for the k return value of the question.

$$\text{Entropy} = H(x_i) = -\frac{1}{\log(n)} \sum_{i=1}^n P(x_i) \log_2 P(x_i)$$

From the wikipedia definition of entropy [https://en.wikipedia.org/wiki/Entropy_\(information_theory\)](https://en.wikipedia.org/wiki/Entropy_(information_theory)).

The $\frac{1}{\log(n)}$ is added to scale branches that have eliminated an example outcome to to have the same max entropy as those that have a larger number of example outcomes remaining.