

## Maximizing information gain – getting the most bang for the buck

In order to split the nodes at the most informative features, we need to define an objective function that we want to optimize via the tree learning algorithm. Here, our objective function is to maximize the information gain at each split, which we define as follows:

$$IG(D_p, f) = I(D_p) - \sum_{j=1}^m \frac{N_j}{N_p} I(D_j)$$

Here,  $f$  is the feature to perform the split,  $D_p$  and  $D_j$  are the dataset of the parent and  $j$ th child node,  $I$  is our impurity measure,  $N_p$  is the total number of samples at the parent node, and  $N_j$  is the number of samples in the  $j$ th child node. As we can see, the information gain is simply the difference between the impurity of the parent node and the sum of the child node impurities – the lower the impurity of the child nodes, the larger the information gain. However, for simplicity and to reduce the combinatorial search space, most libraries (including scikit-learn) implement binary decision trees. This means that each parent node is split into two child nodes,  $D_{left}$  and  $D_{right}$ :

$$IG(D_p, f) = I(D_p) - \frac{N_{left}}{N_p} I(D_{left}) - \frac{N_{right}}{N_p} I(D_{right})$$

Now, the three impurity measures or splitting criteria that are commonly used in binary decision trees are **Gini impurity** ( $I_G$ ), **entropy** ( $I_H$ ), and the **classification error** ( $I_E$ ). Let's start with the definition of entropy for all **non-empty** classes  $p(i|t) \neq 0$ :

$$I_H(t) = - \sum_{i=1}^c p(i|t) \log_2 p(i|t)$$

Here,  $p(i|t)$  is the proportion of the samples that belongs to class  $i$  for a particular node  $t$ . The entropy is therefore 0 if all samples at a node belong to the same class, and the entropy is maximal if we have a uniform class distribution. For example, in a binary class setting, the entropy is 0 if  $p(i=1|t)=1$  or  $p(i=0|t)=0$ . If the classes are distributed uniformly with  $p(i=1|t)=0.5$  and  $p(i=0|t)=0.5$ , the entropy is 1. Therefore, we can say that the entropy criterion attempts to maximize the mutual information in the tree.

Intuitively, the Gini impurity can be understood as a criterion to minimize the probability of misclassification:

$$I_G(t) = \sum_{i=1}^c p(i|t)(1-p(i|t)) = 1 - \sum_{i=1}^c p(i|t)^2$$

Similar to entropy, the Gini impurity is maximal if the classes are perfectly mixed, for example, in a binary class setting ( $c=2$ ):

$$1 - \sum_{i=1}^2 0.5^2 = 0.5$$

However, in practice both the Gini impurity and entropy typically yield very similar results and it is often not worth spending much time on evaluating trees using different impurity criteria rather than experimenting with different pruning cut-offs.

Another impurity measure is the classification error:

$$I_E = 1 - \max \{p(i|t)\}$$