

Decision tree learning

From Wikipedia, the free encyclopedia

Gini impurity

Main article: [Gini coefficient](#)

Used by the CART (classification and regression tree) algorithm, Gini impurity is a measure of how often a randomly chosen element from the set would be incorrectly labeled if it were randomly labeled according to the distribution of labels in the subset. Gini impurity can be computed by summing the probability of each item being chosen times the probability of a mistake in categorizing that item. It reaches its minimum (zero) when all cases in the node fall into a single target category.

To compute Gini impurity for a set of items, suppose i takes on values in $\{1, 2, \dots, m\}$, and let f_i be the fraction of items labeled with value i in the set.

$$I_G(f) = \sum_{i=1}^m f_i(1 - f_i) = \sum_{i=1}^m (f_i - f_i^2) = \sum_{i=1}^m f_i - \sum_{i=1}^m f_i^2 = 1 - \sum_{i=1}^m f_i^2$$

Information gain

Main article: [Information gain in decision trees](#)

Used by the [ID3](#), [C4.5](#) and C5.0 tree-generation algorithms. [Information gain](#) is based on the concept of [entropy](#) from [information theory](#).

$$I_E(f) = - \sum_{i=1}^m f_i \log_2 f_i$$