

Decision Tree

24th November 2021

1 Biến đổi lại gini score, gini index và cách xây dựng decision tree

1.1 Gini score

$$giniscore = 1 - \sum_{i=1}^C p_i^2 = 1 - \sum_{i=1}^C \left(\frac{n_i}{N}\right)^2$$

n_i là số lượng phần tử của lớp i
 N là tổng số phần tử ở node
 $N = \sum_{i=1}^N n_i \rightarrow \sum_{i=1}^N p_i = 1$

Set: $a_i = \frac{n_i}{N}$
 $S = \sum_{i=1}^N a_i^2$
điều kiện: $0 \leq a_i \leq 1$ và $\sum a_i = 1$

1.1.1 S min

We have Bu-nhi-a:

$$(a_1^2 + a_2^2 + \dots + a_N^2) \cdot (1^2 + 1^2 + \dots + 1^2) \geq (a_1 + a_2 + \dots + a_N)^2$$

$$\Leftrightarrow S \cdot N \geq 1$$

$$\Leftrightarrow S \geq \frac{1}{N}$$

$S_{min} = \frac{1}{N}$ khi $a_1 = a_2 = \dots = a_N = \frac{1}{N}$

1.1.2 S max

$$S = \sum_{i=1}^N a_i^2 \leq \left(\sum_{i=1}^N a_i\right)^2 = \sum_{i=1}^N a_i^2 + \sum_{i=1}^N \sum_{j=1}^N a_i \cdot a_j$$

dấu "=" xảy ra khi $a_j = 0$ và $a_{ij} = 0$

1.2 Gini index

$$giniindex = giniparent - \sum_{i=1}^N \frac{n_i}{N} \cdot g_{c_i}$$

1.3 Xây dựng decision tree

Maximize gini-index