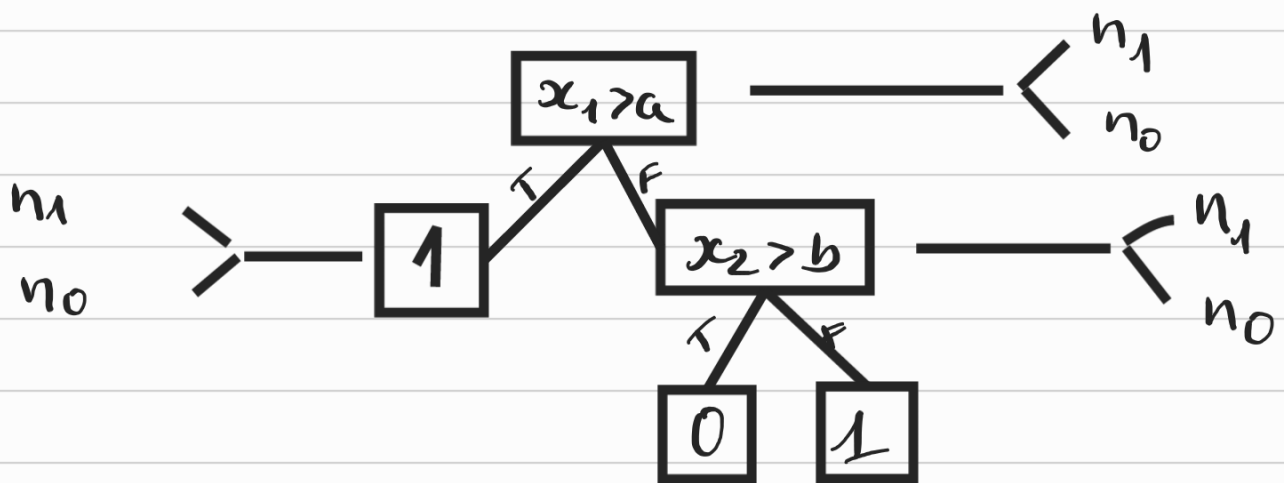


# Gini Score , Gini index

- Để lựa chọn được điều kiện phù hợp nhất cho việc phân tách dữ liệu thành các lớp, máy tính dựa vào các chỉ số để đánh giá. Trong đó có chỉ số gini index.

## ① Xây dựng cây quyết định

VĐ của xây được 1 cây như sau :



$$\text{Gini Score} = 1 - \sum_{i=1}^C (p_i)^2$$

+  $C$  : số lớp cần phân loại

+  $p_i = \frac{n_i}{N}$   $\left\{ \begin{array}{l} n_i : \text{số lượng phần tử lớp thứ } i \\ N : \text{số lượng phần tử ở node} \end{array} \right.$

$$+ \sum_{i=1}^N p_i = 1, \quad N = \sum_{i=1}^N n_i$$

⊕  $\bar{v} \leq \sum p_i \leq 1$  và  $0 \leq p_i \leq 1$  nên:

$$\sum_{i=1}^C (p_i)^2 \leq \left( \sum_{i=1}^C p_i \right)^2 = 1$$

$$\Rightarrow \text{gini score} = 1 - \sum_{i=1}^C (p_i)^2 \geq 0$$

dấu = xảy ra khi  $\exists j \mid p_j = 1$  và  $p_k = 0$   
( $k \neq j$ )

⊕  $\bar{v} \quad C > 1$

$$\sum_{i=1}^C (p_i)^2 \geq \frac{\left( \sum_{i=1}^C p_i \right)^2}{C} = \frac{1}{C} \quad \left( 0 \leq \frac{1}{C} \leq 1 \right)$$

$$\Rightarrow \text{gini score} = 1 - \sum_{i=1}^C (p_i)^2$$

$$\Rightarrow \text{gini score} \leq 1 - \frac{1}{C}$$

Dấu = xảy ra  $\Leftrightarrow p_j = \frac{1}{C} \quad \forall j$

$$\Rightarrow \begin{cases} \max & \text{gini score} = 1 - \frac{1}{C} \\ \min & \text{gini score} = 0 \end{cases}$$

$$\text{Vậy, gini index} = \text{gini}(p) - \sum_{i=1}^K \frac{m_i}{M} \text{gini}(c_i)$$

+  $\text{gini}(p)$  = gini score ở node cha

+  $K$  số node con được tách ra

+  $\text{gini}(c_i)$  chỉ số gini ở node con thứ  $i$

+  $M$  là số phần tử ở node  $p$

+  $m_i$  là số phần tử ở node con thứ  $i$

$$+ \sum_{i=1}^K m_i = M$$

$\Rightarrow$  Để tìm đk tách, ta sẽ thử tất cả thuộc tính,  $\forall$  thuộc tính, thử 1 giá trị chia xem gini index của giá trị chia nào max thì chọn

$\left\{ \begin{array}{l} \text{gini score càng nhỏ càng tốt} \\ \text{gini index càng lớn càng tốt} \end{array} \right.$