

Decision Trees - Parte 2: Índice GINI

$$GINI = 1 - \sum_{i=1}^C p_i^2$$

VarA	VarB	Classe
1	13	1
0	-2	1
0	27	1
1	9	1
0	67	0
0	45	0
1	21	0
0	50	0

D Variável A

Var A = 1:

Classe 1: 2/3

Classe 0: 1/3

$$GINI = 1 - ((2/3)^2 + (1/3)^2) = 0,44$$

Var A = 0:

Classe 1: 2/5

Classe 0: 3/5

$$GINI = 1 - ((2/5)^2 + (3/5)^2) = 0,48$$

$$GINI_{VarA} = \frac{3}{8} \cdot 0,44 + \frac{5}{8} \cdot 0,48 = \underline{\underline{0,46}}$$

D Variável B

Var B < 27:

Classe 1: 3/4

Classe 0: 1/4

$$GINI = 1 - ((3/4)^2 + (1/4)^2) = 0,37$$

__/__/__

S T Q Q S S D

#VONB ≥ 27 :

CLASSE 1: 1/4

GINI = 0,37

CLASSE 0: 3/4

$$GINI_{VONB} = \frac{4}{8} \cdot 0,37 + \frac{4}{8} \cdot 0,37 = 0,37$$

OBS. PARA COMEÇAR A CRIAR NOSSA ÁRVORE DE DECISÃO DEVEMOS ESCOLHER A VARIÁVEL COM O MENOR ÍNDICE GINI

OBS. EM UM CASO HIPOTÉTICO EM QUE OS DADOS ESTÃO BEM DIVIDIDOS (MEIO A MEIO), O ÍNDICE GINI É BAIXO. QUANTO MAIS DIVIDIDO OS DADOS ESTÃO, MENOR O ÍNDICE GINI, ENQUANTO QUE QUANTO MAIS SEPARANDO OS DADOS ESTIVEREM, MAIOR O ÍNDICE GINI, SENDO QUE ESTE ÚLTIMO CASO É O CASO QUE MAIS VAI EXPLICAR OS DADOS.

VonA VonB CLASS

1 13 1
0 -2 1
1 9 1
1 21 0

VonB

≤ 17

≥ 27

VonB

VonB

≤ 13

≥ 13

VonA VonB CLASS

0 27 1
0 07 0
0 45 0
0 50 0