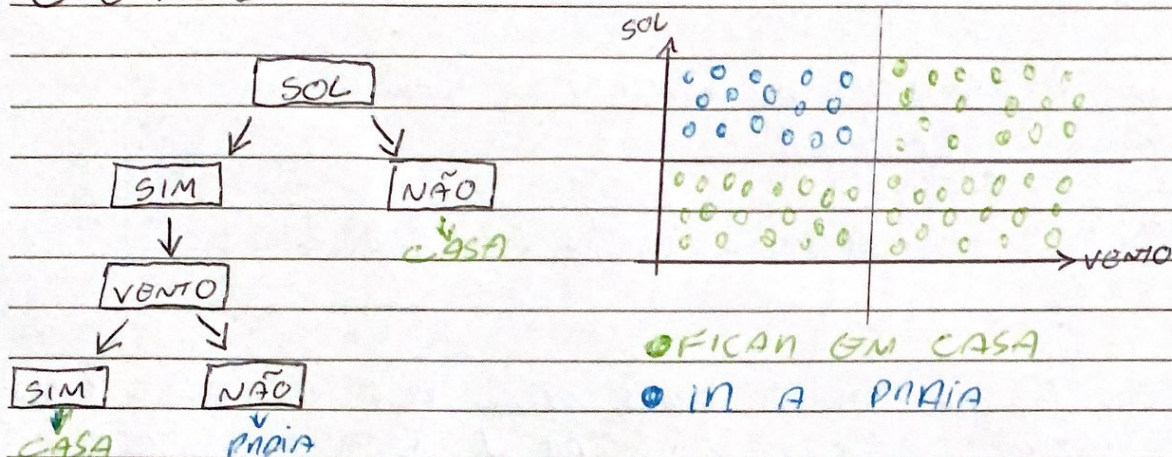
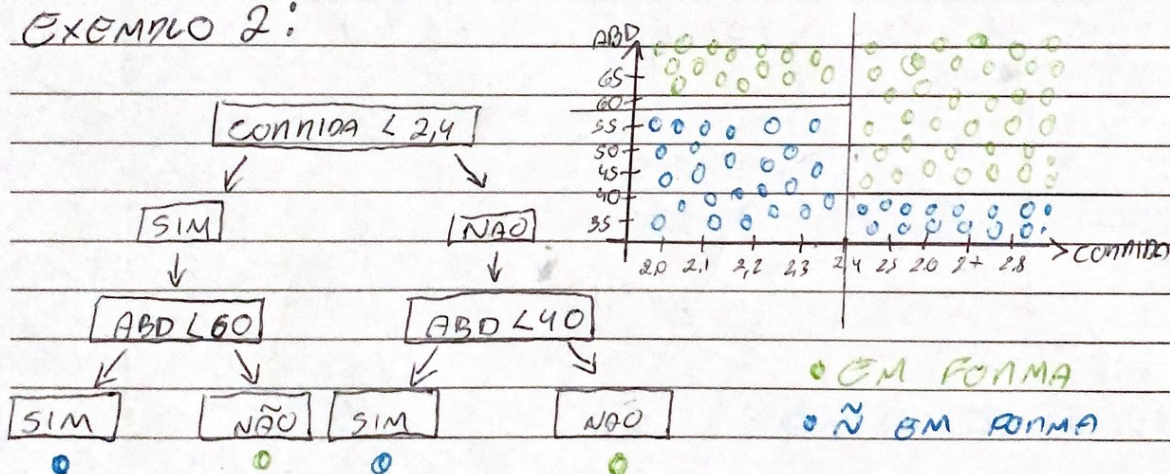


DECISION TREES - PARTE 1: ENTROPIA

EXEMPLO 1:



EXEMPLO 2:



$$\text{GANHO INFORMAÇÃO} = \text{ENTROPIA}_{\text{PAI}} - \sum \text{PESO}_{\text{FILHO}} \cdot \text{ENTROPIA}_{\text{FILHO}}$$

$$\text{ENTROPIA} = -\sum P_i \log_2 P_i, \quad \text{PESO} = \frac{\text{N}^\circ \text{ Amostras Filho}}{\text{N}^\circ \text{ Amostras Pai}}$$

SALONIO	LOCALIZAÇÃO	FUNÇÃO	DECISÃO
ALTO	LONGE	INTERESSANTE	SIM
BAIXO	PERTO	DESINTERESSANTE	NÃO
BAIXO	LONGE	INTERESSANTE	SIM
AUTO	LONGE	DESINTERESSANTE	NÃO
ALTO	PERTO	INTERESSANTE	SIM
BAIXO	LONGE	DESINTERESSANTE	NÃO

___/___/___

S T Q Q S S D

$$P_S = 3/6 = 1/2$$

$$P_N = 3/6 = 1/2$$

$$\text{Entropia}_D = -\left(\frac{1}{2} \log_2 \frac{1}{2} + \frac{1}{2} \log_2 \frac{1}{2}\right)$$

$$= 1$$

Caso a proporção de SIM para NÃO fosse de 5 para 1, teríamos:

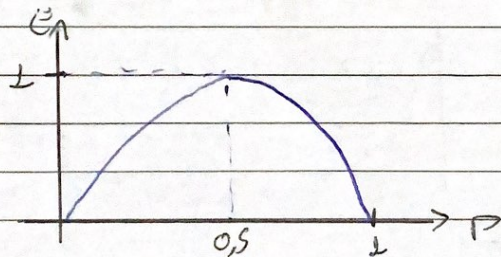
$$P_S = 5/6$$

$$\text{Entropia}_D = -\left(\frac{5}{6} \log_2 \frac{5}{6} + \frac{1}{6} \log_2 \frac{1}{6}\right)$$

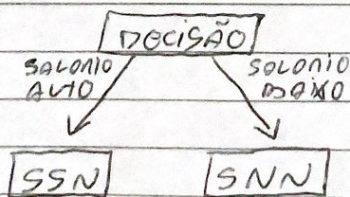
$$P_N = 1/6$$

$$= 0,65$$

OBS. Quanto mais os dados estão bem divididos (bem na metade), mais próximo de 1 será a entropia, enquanto que mais dispersos ou desbalanceados, mais próximo de zero será a entropia



#VARIÁVEL SALONIO



▷ RAMO/FILHO SALONIO ALTO:

$$P_S = 2/3$$

$$P_N = 1/3$$

$$E = -\left(\frac{2}{3} \log_2 \frac{2}{3} + \frac{1}{3} \log_2 \frac{1}{3}\right)$$

$$= 0,92$$

$$P_{\text{DSO}} = 3/6 = 1/2$$

$$G_1 = 1 - \left(\frac{1}{2} \cdot 0,92 + \frac{1}{2} \cdot 0,92\right)$$

$$= 0,08$$

▷ RAMO/FILHO SALONIO BAIXO:

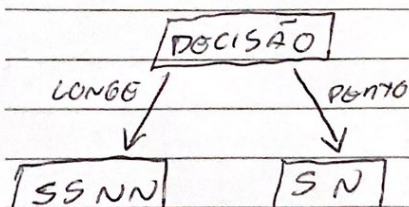
$$P_S = 1/3$$

$$E = 0,92$$

$$P_N = 2/3$$

$$P_{\text{DSO}} = 3/6 = 1/2$$

VARIÁVEL LOCALIZAÇÃO



▷ RAMO LONGE

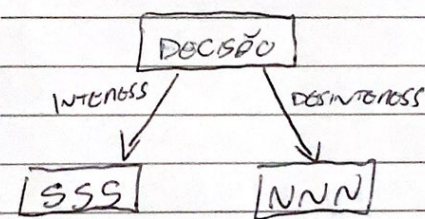
$$\begin{aligned}
 & \bullet P_S = 2/4 = 1/2 \quad \bullet P_{CSO} = 4/6 \\
 & \bullet P_N = 2/4 = 1/2 \quad \bullet E = 1
 \end{aligned}$$

▷ RAMO PERTO

$$\begin{aligned}
 & \bullet P_S = 1/2 \quad \bullet P_{CSO} = 2/6 \\
 & \bullet P_N = 1/2 \quad \bullet E = 1
 \end{aligned}$$

$$\begin{aligned}
 GI &= 1 - (4/6 \cdot 1 + 2/6 \cdot 1) \\
 &= 0
 \end{aligned}$$

VARIÁVEL FUNÇÃO



▷ RAMO INTERESSANTE

$$\bullet E = 0$$

▷ RAMO DESINTERESSANTE

$$\bullet E = 0$$

$$GI = 1$$

OBS. A variável que tiver o maior ganho de informação será a primeira variável a ser utilizada para a criação do árvore de decisão.

OBS. No caso específico da variável função, como o GI foi 1, essa árvore já seria nossa árvore final.

OBS. Depois de escolher a variável com maior GI, a gente continua construindo nosso árvore de decisão como fizemos até aqui porém considerando esses nós de parada como sendo os próximos pontos de parada para calcular o restante da árvore.