

CÁLCULO DA ENTROPIA

Em nossa aula sobre árvore de decisão, usamos a seguinte base de dados como exemplo para construção da árvore:

Aspecto	Temperatura	Umidade	Vento	Jogar
Sol	Quente	Alta	Não	Não
Sol	Quente	Alta	Sim	Não
Nublado	Quente	Alta	Não	Sim
Chuva	Branda	Alta	Não	Sim
Chuva	Fria	Normal	Não	Sim
Chuva	Fria	Normal	Sim	Não
Nublado	Fria	Normal	Sim	Sim
Sol	Branda	Alta	Não	Não
Sol	Fria	Normal	Não	Sim
Chuva	Branda	Normal	Não	Sim
Sol	Branda	Normal	Sim	Sim
Nublado	Branda	Alta	Sim	Sim
Nublado	Quente	Normal	Não	Sim
Chuva	Branda	Alta	Sim	Não

O cálculo da entropia é obtido da seguinte forma:

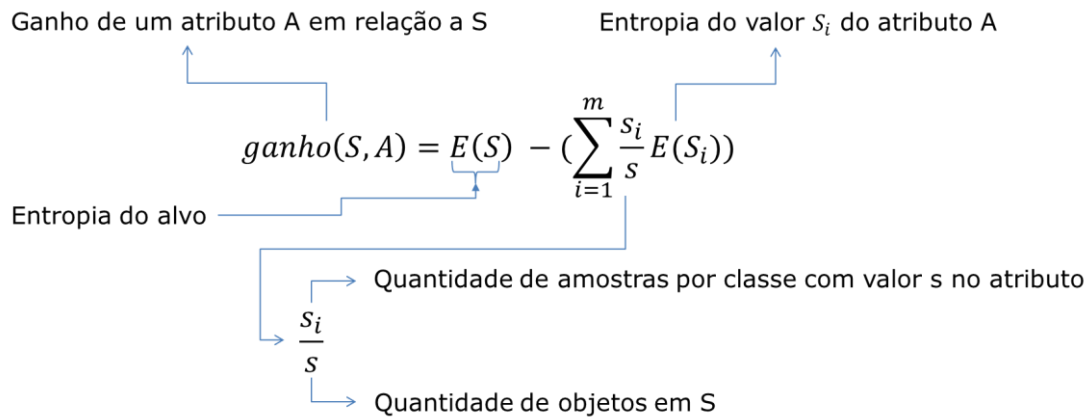
Amostras de um atributo Quantidade de classes da base

$$E(S) = - \sum_{i=1}^m p_i \log_2 p_i$$

p_i é a probabilidade de que uma amostra pertença a uma classe c_i

$$p_i = \frac{s_i}{s}$$

IMPACTA



Vale lembrar que os termos amostra, objeto ou ainda instância são equivalentes para se referir a uma linha de uma base de dados.

Então, a entropia da nossa resposta esperada (alvo) é:

$$S = 14 \text{ amostras}$$

$$S_{sim} = 9$$

$$S_{n\tilde{a}o} = 5$$

$$E(9,5) = -\left(\frac{9}{14} \log_2 \frac{9}{14} + \frac{5}{14} \log_2 \frac{5}{14}\right)$$

$$E(9,5) = 0.940 \text{ bits}$$

Aspecto (Sol, Nublado, Ensolarado)

Sol ocorre 5 vezes na base em que 2 vezes a resposta esperada é **jogar** e 3 **não jogar**:

$$Sol (S_{sim} = 2, S_{n\tilde{a}o} = 3, S = 5)$$

Nublado ocorre 4 vezes em todas as ocorrências a resposta é **não jogar**:

$$Nublado (S_{sim} = 4, S_{n\tilde{a}o} = 0, S = 4)$$

Chuvoso ocorre 5 vezes na base em que 3 vezes a resposta é **jogar** e 2 **não jogar**:

IMPACTA

Chuva ($S_{sim} = 3, S_{n\tilde{a}o} = 2, S = 5$)

Agora, calculamos a entropia de cada um dos valores do aspecto:

$$E(S) = - \sum_{i=1}^m p_i \log_2 p_i$$

$$Sol = -\left(\frac{2}{5} \log_2 \frac{2}{5} + \frac{3}{5} \log_2 \frac{3}{5}\right) = 0.971$$

$$Nublado = -\left(\frac{4}{4} \log_2 \frac{4}{4} + \frac{0}{4} \log_2 \frac{0}{4}\right) = 0$$

$$Chuva = -\left(\frac{3}{5} \log_2 \frac{3}{5} + \frac{2}{5} \log_2 \frac{2}{5}\right) = 0.971$$

Temos, então, que o ganho de informação fornecido pelo atributo **Aspecto** é:

$$ganho(S, Aspecto) = 0.940 - \left(0.971 * \frac{5}{14} + 0 * \frac{4}{14} + 0.971 * \frac{5}{14}\right) = 0.246$$

Fazendo o cálculo para os demais atributos obtemos:

$$ganho(S, Temperatura) = 0.940 - \left(1 * \frac{4}{14} + 0.918 * \frac{6}{14} + 0.811 * \frac{4}{14}\right) = 0.029$$

$$ganho(S, Umidade) = 0.940 - \left(0.985 * \frac{7}{14} + 0.592 * \frac{7}{14}\right) = 0.152$$

$$ganho(S, Vento) = 0.940 - \left(0.811 * \frac{8}{14} + 1 * \frac{6}{14}\right) = 0.048$$

Nosso objetivo é encontrar o atributo com maior ganho de informação, por isso o escolhido é o atributo **Aspecto**.