

2.

(Calcularea câștigului de informație pe "decision stumps")

Studentul Timmy dorește să știe cum [ar trebui să procedeze cel mai bine ca] să promoveze examenul de învățare automată. Pentru aceasta, a cules informații de la studenții care au urmat acest curs în anii precedenți și apoi a decis să-și construiască un *model* bazat pe arbori de decizie. A a colectat în total nouă *instanțe*/exemplare, descrise cu ajutorul a două *trăsături* (văzute mai departe ca două variabile aleatoare,  $S$  și  $A$ ): „este bine să stai și să înveți până noaptea târziu înainte de examen” ( $S$ ) și „este bine să mergi la toate cursurile și seminariile” ( $A$ ). Timmy dispune acum de următoarele „statistici” (care sunt de fapt *partiționări* ale datelor sale):

$$\begin{aligned} \text{Set}(\text{all}) &= \{5+, 4-\} \\ \text{Set}(S+) &= \{3+, 2-\}, \text{Set}(S-) = \{2+, 2-\} \\ \text{Set}(A+) &= \{5+, 1-\}, \text{Set}(A-) = \{0+, 3-\} \end{aligned}$$

Presupunând că se folosește drept criteriu de selecție a celei mai bune trăsături câștigul maxim de informație, ce trăsătură va alege Timmy? Care este valoarea câștigului de informație?

Puteți folosi la calcule următoarele aproximații:

$N$	1	2	3	4	5	6	7	8	9
$\log_2 N$	0	1	1.58	2	2.32	2.58	2.81	3	3.17

Calculăm mai întâi entropia variabilei  $A$ :

$$H(A) = -\frac{1}{2} p(x) \cdot \log p(x) \cdot$$

$$\begin{aligned} H(A) &= H(5, 4) = \log_2 5 \cdot \log \frac{5}{9} + \frac{4}{9} \cdot \log \frac{4}{9} = \\ &= \frac{9 \log 9 - 4 \log 4 - 5 \log 5}{9} = \frac{28.52932 - 8 - 11.6096}{9} = \end{aligned}$$

$$\approx \frac{0.99104}{9}$$

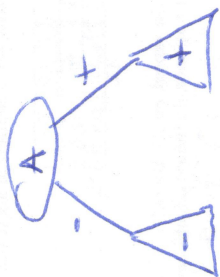
$$\begin{aligned} H(S) &= \frac{1}{2} \cdot \log_2 2 + \frac{1}{2} \cdot \log_2 2 = 1 \\ H(S+) &= \frac{1}{2} \cdot \log_2 2 + \frac{1}{2} \cdot \log_2 2 = 1 \\ H(S-) &= \frac{1}{2} \cdot \log_2 2 + \frac{1}{2} \cdot \log_2 2 = 1 \end{aligned}$$

$$IG(S) = H(A) - H(S) = 0.99104 - 1 = -0.00896$$

$$\begin{aligned} H_0(A) &= \frac{2}{9} \cdot H(0, 3) + \frac{6}{9} \cdot H(1, 5) = \\ &= \frac{2}{9} \left( \frac{1}{6} \cdot \log 6 + \frac{5}{6} \cdot \log \frac{6}{5} \right) = \frac{6 \log 6 - 5 \log 5}{9} = \\ &= \frac{15.509475 - 11.6096}{9} = \frac{0.433}{9} \Rightarrow \\ IG(0, A) &= H(\text{all}) - H_0(A) \Rightarrow IG(0, A) = \frac{0.55472}{9} \end{aligned}$$

Care  $IG(0, A) > IG(1, S) \Rightarrow$

Arborele va avea astfel:



foarte bine!

Concluzia: pt. a promova la examen, este bine să mergi la toate cursurile și seminariile.