

### 3.1 Arbori de decizie — Probleme rezolvate

1. (Arbori de decizie;  
optimalitate, relativ la numărul de noduri)

Reprezentați arborele/arborii de decizie care are/au numărul minim de noduri posibile și corespunde/corespund funcției boolene  $(\neg A \vee B) \wedge \neg(C \wedge A)$  definită peste atributele boolene  $A, B$  și  $C$ .

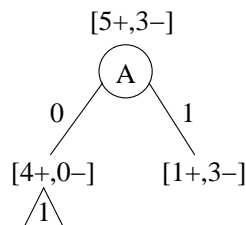
Răspuns:

Vom determina arborele de decizie optimal (ca număr de noduri) parcurgând în mod *exhaustiv spațiul de versiuni*, adică mulțimea tuturor arborilor de decizie (construiți cu variabilele  $A, B$  și  $C$ ) care sunt *consistenți* cu funcția dată. Așadar, vom examina ce se întâmplă când în nodul rădăcină se pun pe rând atributele  $A, B$  și respectiv  $C$ .

Notăm cu  $X$  funcția  $(\neg A \vee B) \wedge \neg(C \wedge A)$ , ale cărei valori sunt date în tabelul alăturat.

$A$	$B$	$C$	$X$
0	0	0	1
0	0	1	1
0	1	0	1
0	1	1	1
1	0	0	0
1	0	1	0
1	1	0	1
1	1	1	0

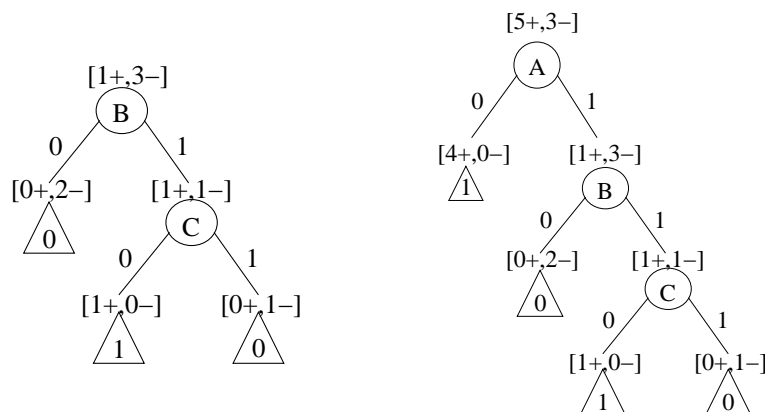
- *Cazul 1:* Dacă în nodul rădăcină se plasează atributul  $A$ , va rezulta următoarea re-partiționare a mulțimii de exemple pentru conceptul  $X$ :



Subarborele drept va trebui să reprezinte arborele de decizie pentru funcția

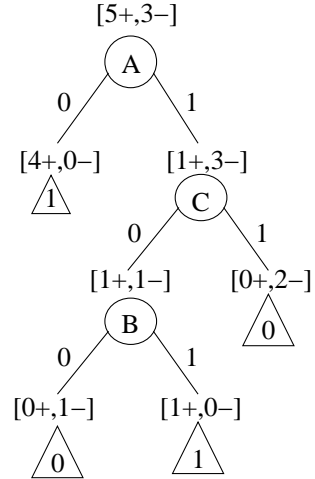
$$X_1 = X[A/1] = (\neg 1 \vee B) \wedge \neg(C \wedge 1) = B \wedge \neg C,$$

pentru care o reprezentare optimă este redată mai jos, în partea stângă:

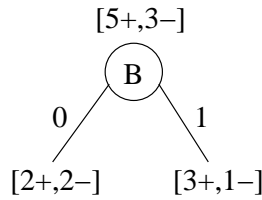


Prin urmare, un arbore optim (ca număr de noduri) care are variabila  $A$  în nodul rădăcină este cel reprezentat mai sus, în partea dreaptă.

*Observație:* Evident, există încă un arbore optim care are variabila  $A$  în nodul rădăcină (el corespunde unei alte reprezentări optimale a conjuncției  $B \wedge \neg C$  față de cea de mai sus). Vedeți desenul alăturat.



- *Cazul 2:* Dacă în nodul rădăcină se alege atributul  $B$ ,



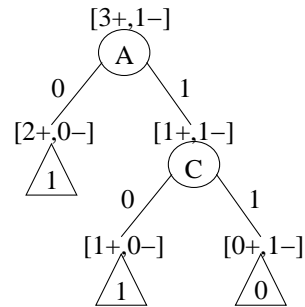
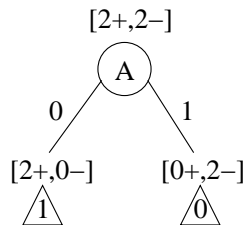
subarborele stâng și subarborele drept trebuie să reprezinte arborele de decizie pentru funcțiile

$$X_2 = X[B/0] = (\neg A \vee 0) \wedge \neg(C \wedge A) = \neg A \wedge (\neg C \vee \neg A) = (\neg A \wedge \neg C) \vee \neg A = \neg A,$$

și respectiv

$$X_3 = X[B/1] = (\neg A \vee 1) \wedge \neg(C \wedge A) = 1 \wedge (\neg C \vee \neg A) = \neg C \vee \neg A,$$

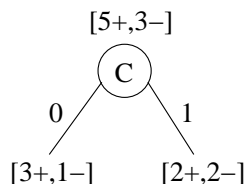
care au ca reprezentări optime arborii de mai jos:



Pentru arborele din dreapta există un arbore de decizie echivalent, obținut prin interschimbarea lui  $A$  cu  $C$ .

Prin urmare, orice arbore optim având variabila  $B$  în rădăcină are 3 niveluri și 4 noduri, așadar cu un nod (de test) mai mult decât cel determinat în primul caz.

- *Cazul 3:* În sfârșit, când în nodul rădăcină se alege atributul  $C$ ,



subarborii drept și stâng trebuie să reprezinte arborele de decizie pentru funcțiile

$$X_4 = X[C/0] = (\neg A \vee B) \wedge \neg(0 \wedge A) = (\neg A \vee B) \wedge \neg 0 = \neg A \vee B,$$

și respectiv

$$X_5 = X[C/1] = (\neg A \vee B) \wedge \neg(1 \wedge A) = (\neg A \vee B) \wedge \neg A = \neg A$$

Urmând un raționament similar cu cel de la cazul anterior, putem spune că orice arbore optim cu atributul  $C$  în rădăcină are 3 niveluri și 4 noduri, cu un nod (de test) mai mult decât cel determinat în primul caz.

Așadar, putem concluziona că arborii de decizie optimi corespunzători funcției date sunt cei determinați în primul caz.

2.

(Algoritmul ID3: aplicare)

■● *CMU, 2002 spring, A. Moore, midterm example questions, pr. 2*

Ai naufragiat pe o insulă pustie, unde nu găsești niciun alt fel de hrană decât ciuperci. Despre unele dintre aceste ciuperci se știe că sunt otrăvitoare, despre altele se știe că sunt comestibile, iar despre restul nu se știe ce fel sunt. Ai rămas singur pe insulă — foștii tăi camarazi, fiind epuizați de foame, au folosit metoda ‘trial and error’... — și ai la dispoziție următoarele date:

Exemplu	<i>Ușoară</i>	<i>Mirositoare</i>	<i>ArePete</i>	<i>Netedă</i>	<i>Comestibilă</i>
<i>A</i>	1	0	0	0	1
<i>B</i>	1	0	1	0	1
<i>C</i>	0	1	0	1	1
<i>D</i>	0	0	0	1	0
<i>E</i>	1	1	1	0	0
<i>F</i>	1	0	1	1	0
<i>G</i>	1	0	0	1	0
<i>H</i>	0	1	0	0	0
<i>U</i>	0	1	1	1	?
<i>V</i>	1	1	0	1	?
<i>W</i>	1	1	0	0	?

Atunci când nu vei mai avea la dispoziție pentru a supraviețui decât ciuperci  $U$ ,  $V$ , sau  $W$ , ai putea estima care dintre ele sunt comestibile, folosind arbori de decizie.

În primele trei întrebări care urmează, ne vom referi la ciupercile  $A - H$ :

- Care este entropia atributului *Comestibilă*?
- Doar privind datele — adică fără a face explicit calculul câștigului de informație (engl., information gain) pentru cele patru atribute — poți determina ce atribut vei alege ca rădăcină a arborelui de decizie?
- Calculează câștigul de informație pentru atributul pe care l-ai ales la întrebarea precedentă.
- Elaborează întregul arbore de decizie ID3 bazat pe datele din tabel și apoi clasifică ciupercile U, V, W.
- Exprimă cu ajutorul calculului propozițional (logica predicatelor de ordinul 0) clasificarea produsă de arborele de decizie obținut. (*Comestibilă*  $\leftrightarrow \dots$ )
- Există vreun risc dacă vei consuma ciuperci care au fost clasificate de arborele de decizie ca fiind comestibile? De ce da? sau, de ce nu?

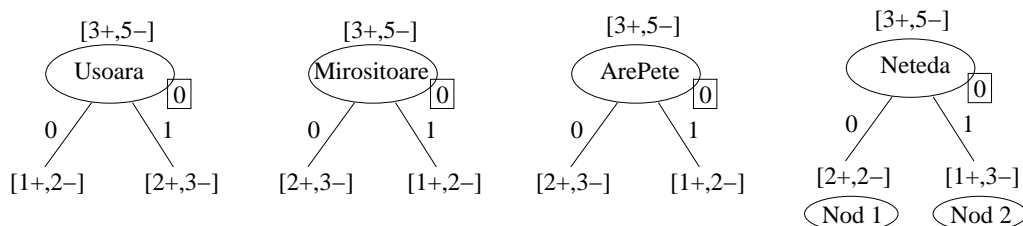
**Răspuns:**

- a. Entropia atributului *Comestibilă* este:

$$\begin{aligned} H_{Comestibilă} &= H[3+, 5-] \stackrel{def.}{=} -\frac{3}{8} \log_2 \frac{3}{8} - \frac{5}{8} \log_2 \frac{5}{8} = \frac{3}{8} \log_2 \frac{8}{3} + \frac{5}{8} \log_2 \frac{8}{5} = \\ &= \frac{3}{8} 3 - \frac{3}{8} \log_2 3 + \frac{5}{8} 3 - \frac{5}{8} \log_2 5 = 3 - \frac{3}{8} \log_2 3 - \frac{5}{8} \log_2 5 \approx \\ &\approx 0.9544 \end{aligned}$$

Notăția  $[3+, 5-]$  simbolizează o mulțime partiționată în 3 exemple pozitive și 5 exemple negative. Vom folosi acest gen de notație peste tot în continuare, cu mici variații determinate de valorile pe care le poate lua atributul de ieșire. De exemplu, dacă vorbim despre o mulțime cu 5 obiecte roșii, 3 albastre și 4 verzi, am putea nota:  $[5R, 3A, 4V]$ .

- b. În rădăcina arborelui de decizie se alege atributul care aduce cel mai mare câștig de informație. Adică, atributul care, intuitiv vorbind, partiționează cel mai bine datele de antrenament în raport cu atributul de ieșire. În cazul nostru, variantele pe care le avem la dispoziție pentru rădăcina arborelui (nodul 0) sunt:



Este ușor de observat că atributele *Ușoară*, *Mirositoare* și *ArePete* împart mulțimea exemplilor în mod similar: o submulțime cu 3 elemente, dintre care unul este pozitiv iar două sunt negative, și o submulțime cu 5 elemente, dintre care două sunt pozitive, iar trei sunt negative.

Dacă am considera un arbore de decizie cu un singur nod de test în care plasăm atributul *Netedă*, atunci numărul minim de erori la antrenare pe care îl putem obține este 3, utilizând următoarea clasificare:

- $Neted\grave{a} = 0 : Comestibil\grave{a} = 1 \Rightarrow$  ciupercile  $E$  și  $H$  sunt clasificate greșit
- $Neted\grave{a} = 1 : Comestibil\grave{a} = 0 \Rightarrow$  ciuperca  $C$  este clasificată greșit

Dacă, în schimb, vom pune în rădăcina arborelui de decizie unul dintre celelalte trei atribute, spre exemplu atributul  $Ușoar\grave{a}$ , și dacă vom lua votul majoritar în fiecare nod descendent din nodul rădăcină, eroarea rezultată la antrenare va fi aceeași ca mai sus ( $3/8$ ), însă toate instanțele vor fi clasificate la fel (și anume, negativ). Dacă nu lucrăm cu vot majoritar pentru ambii descendenți, ci doar pentru cel cu entropie mai mică (în vreme ce pentru celălalt nod descendent luăm decizia contrară), se observă că pentru atributul  $Ușoar\grave{a}$  vom obține 4 erori pe setul de antrenament, iar pentru atributul  $Neted\grave{a}$  vom obține 3 erori.

Sumarizând, suntem înclinați să credem că ar fi o alegere sensibil mai bună să punem în rădăcină atributul  $Neted\grave{a}$ . Pentru o justificare numerică riguroasă a acestei alegeri folosind criteriul maximizării câștigului de informație, vezi punctul  $d$ .

c. Pentru a obține câștigul de informație pentru atributul  $Neted\grave{a}$ , se fac calculele:

$$\begin{aligned}
 H_{0/Neted\grave{a}} &\stackrel{def.}{=} \frac{4}{8}H[2+, 2-] + \frac{4}{8}H[1+, 3-] = \frac{1}{2} \cdot 1 + \frac{1}{2} \left( \frac{1}{4} \log_2 \frac{4}{1} + \frac{3}{4} \log_2 \frac{4}{3} \right) \\
 &= \frac{1}{2} + \frac{1}{2} \left( \frac{1}{4} \cdot 2 + \frac{3}{4} \cdot 2 - \frac{3}{4} \log_2 3 \right) = \frac{1}{2} + \frac{1}{2} \left( 2 - \frac{3}{4} \log_2 3 \right) \\
 &= \frac{1}{2} + 1 - \frac{3}{8} \log_2 3 = \frac{3}{2} - \frac{3}{8} \log_2 3 \approx 0.9056 \\
 IG_{0/Neted\grave{a}} &\stackrel{def.}{=} H_{Comestibil\grave{a}} - H_{0/Neted\grave{a}} \\
 &= 0.9544 - 0.9056 = 0.0488
 \end{aligned}$$

În cele de mai sus am notat cu  $H_{0/Neted\grave{a}}$  entropia *partiției* [mulțimii de exemple de antrenament] determinate de alegerea atributului  $Neted\grave{a}$  în nodul 0,<sup>56</sup> iar cu  $IG_{0/Neted\grave{a}}$  câștigul de informație corespunzător acestei alegeri. În general, prin notația  $H_{n/A}$  vom înțelege entropia partiției determinate de alegerea atributului  $A$  în nodul  $n$ .

d. Arborele de decizie ID3 se construiește pornind din rădăcină și alegând atributul pentru fiecare nod de test în modul următor:

**Nodul 0 (rădăcina):**

Să verificăm dacă alegerea făcută la punctul  $b$  este cea corectă:

$$\begin{aligned}
 H_{0/Ușoar\grave{a}} &\stackrel{def.}{=} \frac{3}{8}H[1+, 2-] + \frac{5}{8}H[2+, 3-] \\
 &= \frac{3}{8} \left( \frac{1}{3} \log_2 \frac{3}{1} + \frac{2}{3} \log_2 \frac{3}{2} \right) + \frac{5}{8} \left( \frac{2}{5} \log_2 \frac{5}{2} + \frac{3}{5} \log_2 \frac{5}{3} \right) \\
 &= \frac{3}{8} \left( \frac{1}{3} \log_2 3 + \frac{2}{3} \log_2 3 - \frac{2}{3} \cdot 1 \right) + \frac{5}{8} \left( \frac{2}{5} \log_2 5 - \frac{2}{5} \cdot 1 + \frac{3}{5} \log_2 5 - \frac{3}{5} \log_2 3 \right) \\
 &= \frac{3}{8} \left( \log_2 3 - \frac{2}{3} \right) + \frac{5}{8} \left( \log_2 5 - \frac{3}{5} \log_2 3 - \frac{2}{5} \right)
 \end{aligned}$$

<sup>56</sup>Mai riguros, folosind terminologia din *Teoria informației*, vom spune că notația  $H_{0/Neted\grave{a}}$  se referă la *entropia condițională medie* a atributului de ieșire *Comestibilă* în raport cu atributul de intrare *Netedă*.

$$\begin{aligned}
&= \frac{3}{8} \log_2 3 - \frac{2}{8} + \frac{5}{8} \log_2 5 - \frac{3}{8} \log_2 3 - \frac{2}{8} \\
&= \frac{5}{8} \log_2 5 - \frac{4}{8} \approx 0.9512
\end{aligned}$$

Urmează că

$$IG_{0/U\text{șoară}} \stackrel{\text{def.}}{=} H_{Comestibilă} - H_{0/U\text{șoară}} = 0.9544 - 0.9512 = 0.0032,$$

deci

$$IG_{0/U\text{șoară}} = IG_{0/Mirositoare} = IG_{0/ArPete} = 0.0032 < IG_{0/Netedă} = 0.0488$$

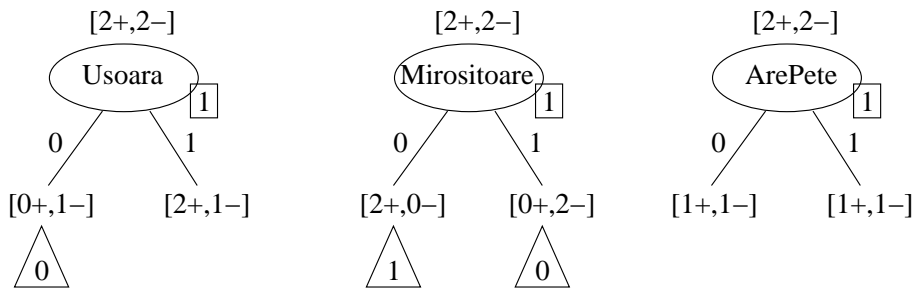
Am avut deci dreptate să alegem atributul *Netedă* la punctul b.

**Observație importantă:** În loc să fi calculat efectiv aceste câștiguri de informație, pentru a determina atributul cel mai „bun“, ar fi fost suficient să elaborăm un *raționament* de tip *relațional*, bazat pe comparația dintre valorile entropiilor condiționale medii  $H_{0/Netedă}$  și  $H_{0/U\text{șoară}}$ :

$$\begin{aligned}
IG_{0/Netedă} > IG_{0/U\text{șoară}} &\Leftrightarrow H_{0/Netedă} < H_{0/U\text{șoară}} \\
&\Leftrightarrow \frac{3}{2} - \frac{3}{8} \log_2 3 < \frac{5}{8} \log_2 5 - \frac{1}{2} \Leftrightarrow 12 - 3 \log_2 3 < 5 \log_2 5 - 4 \\
&\Leftrightarrow 16 < 5 \log_2 5 + 3 \log_2 3 \Leftrightarrow 16 < 11.6096 + 4.7548 \text{ (adev.)}
\end{aligned}$$

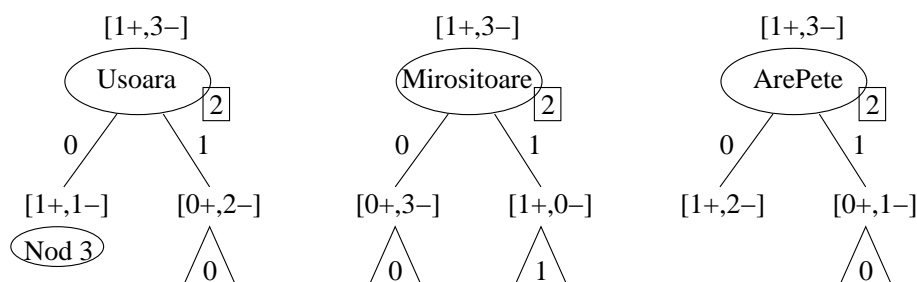
Așa vom proceda adeseori în rezolvările de la acest capitol.

**Nodul 1:** Trebuie să clasificăm acele exemple care au  $Netedă = 0$ ; avem de ales între 3 atribute - *Ușoară*, *Mirositoare* și *ArPete*.



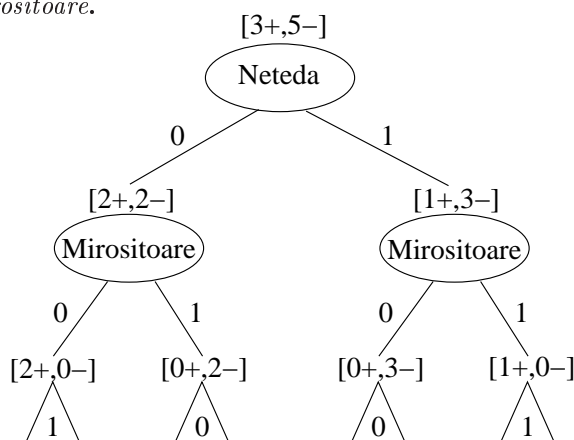
Avem  $H_{1/Mirositoare} = \frac{2}{4} H[2+, 0-] + \frac{2}{4} H[0+, 2-] = 0$ . Oricare ar fi valorile pentru  $H_{1/ArPete}$  și  $H_{1/U\text{șoară}}$ , întrucât știm că entropia are întotdeauna valori nenegative, rezultă că atributul *Mirositoare* maximizează în nodul 1 câștigul de informație. În imaginea de mai sus valorile din triunghi reprezintă decizia luată de subarboarele construit în nodul-frunză respectiv.

**Nodul 2:** Avem de clasificat exemplele pentru care  $Netedă = 1$ . Atributele disponibile sunt: *Ușoară*, *Mirositoare* și *ArPete*.



Evident,  $H_{2/Mirositoare} = \frac{3}{4}H[0+,3-] + \frac{1}{4}H[1+,0-] = \frac{3}{4} \cdot 0 + \frac{1}{4} \cdot 0 = 0$  Așadar, pentru nodul 2 putem alege atributul *Mirositoare*.

Arborele complet arată astfel:



Parcurgând arborele construit, ciupercile  $U$ ,  $V$  și  $W$  vor fi clasificate astfel:

$U$	$Netedă = 1, Mirositoare = 1 \Rightarrow Comestibilă = 1$
$V$	$Netedă = 1, Mirositoare = 1 \Rightarrow Comestibilă = 1$
$W$	$Netedă = 0, Mirositoare = 1 \Rightarrow Comestibilă = 0$

e.  $Comestibilă \leftrightarrow (\neg Netedă \wedge \neg Mirositoare) \vee (Netedă \wedge Mirositoare)$

Același lucru poate fi exprimat și sub forma unui pseudo-cod *if ... then Comestibilă else  $\neg$  Comestibilă*:

```

IF      (Netedă = 0 AND Mirositoare = 0) OR
        (Netedă = 1 AND Mirositoare = 1)
THEN   Comestibilă;
ELSE    $\neg$ Comestibilă;

```

f. Arborele de decizie produs de către algoritmul ID3 elaborat mai sus este consistent cu datele de antrenament pe care le-am avut la dispoziție (fiindcă aceste date sunt necontradictorii). Întrucât în realitate clasificarea poate depinde și de alte trăsături/informații decât cele de care dispunem noi, nu avem garanția că arborele ID3 face identificarea corectă a etichetei/clasei pentru toate instanțele din setul de test. Așadar, nu putem fi siguri că nu ne vom îmbolnăvi dacă vom consuma ciupercile  $U$  și  $V$ , sau că ne vom îmbolnăvi dacă vom consuma ciuperca  $W$ . În multe aplicații practice, calitatea unui model de învățare automată (în cazul de față, un arbore de decizie) se verifică pe un set de *date de validare*.

3. (Algoritmul ID3, aplicat pe expresii booleene;  
exploatarea simetriilor operațiilor  $\vee, \wedge$  în alegerea nodurilor;  
analiza „optimalității” arborelui ID3, ca număr de noduri de test)

\* prelucrare de Liviu Ciortuz după  
Tom Mitchell, “Machine Learning”, 1997, ex. 3.1.b

Considerăm următoarea funcție booleană:  $A \vee (B \wedge C)$ . Presupunem că această funcție este deja definită — adică valoarea ei este cea cunoscută din logica propozițiilor —, însă dorim să o reprezentăm ca arbore de decizie.

a. Aplicați algoritmul ID3 [tabelei de adevăr corespunzătoare] acestei funcții.  
*Observație:* Dacă exploatați simetriile, veți avea nevoie doar de puține calcule, altfel vă veți complica inutil.

b. Arborele ID3 obținut la punctul precedent este optimal?

Alfel spus, puteți găsi un alt arbore de decizie, de adâncime mai mică sau cu număr mai mic de noduri (comparativ cu arborele obținut la punctul a), care să reprezinte această funcție? (Țineți cont că în fiecare nod al unui arbore de decizie se poate testa un singur atribut.)

Răspuns:

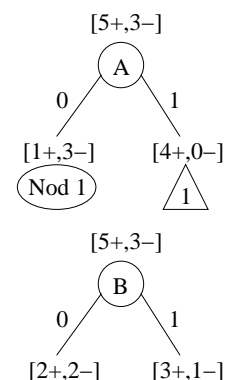
a. Observăm că funcția dată este simetrică în  $B$  și  $C$ , datorită comutativității funcției logice  $\wedge$ . O consecință a acestui fapt este că dacă, pe parcursul algoritmului ID3, avem de ales (și) între cele două atribute este nevoie să-l studiem doar pe unul dintre ele, celălalt comportându-se identic.

$A$	$B$	$C$	$Y = A \vee (B \wedge C)$
0	0	0	0
0	0	1	0
0	1	0	0
0	1	1	1
1	0	0	1
1	0	1	1
1	1	0	1
1	1	1	1

*Nodul 0 (rădăcină):*

$$\begin{aligned}
 H_{0/A} &= \frac{4}{8}H[1+, 3-] + \frac{4}{8}H[4+, 0-] = \\
 &= \frac{1}{2}H[1+, 3-] + \frac{1}{2} \cdot 0 = \\
 &= \frac{1}{2}H[1+, 3-]
 \end{aligned}$$

$$\begin{aligned}
 H_{0/B} &= \frac{4}{8}H[2+, 2-] + \frac{4}{8}H[3+, 1-] = \\
 &= \frac{1}{2} \cdot 1 + \frac{1}{2}H[3+, 1-] = \\
 &= \frac{1}{2} + \frac{1}{2}H[1+, 3-]
 \end{aligned}$$



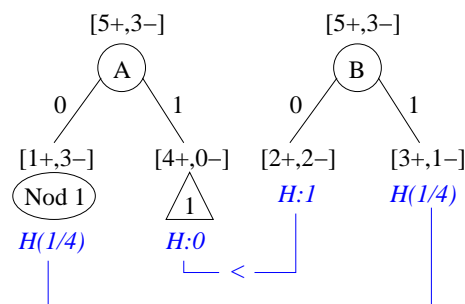
Este evident că  $H_{0/A} < H_{0/B}$ , deci vom alege atributul  $A$  în rădăcină.

*Observație importantă:*

La aceeași concluzie se putea ajunge *imediat* pe baza unui *raționament calitativ*, și anume, comparând atent cei doi arbori („compași de decizie”) de mai

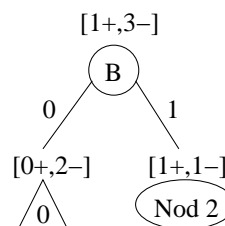


sus. Mai precis, vom compara (două câte două) entropiile condiționale specifice din nodurile descendente, precum și ponderile cu care se combină aceste entropii în scrierea entropiilor condiționale medii corespunzătoare atributelor  $A$  și  $B$ . Putem pune în evidență acest fapt în figura următoare,



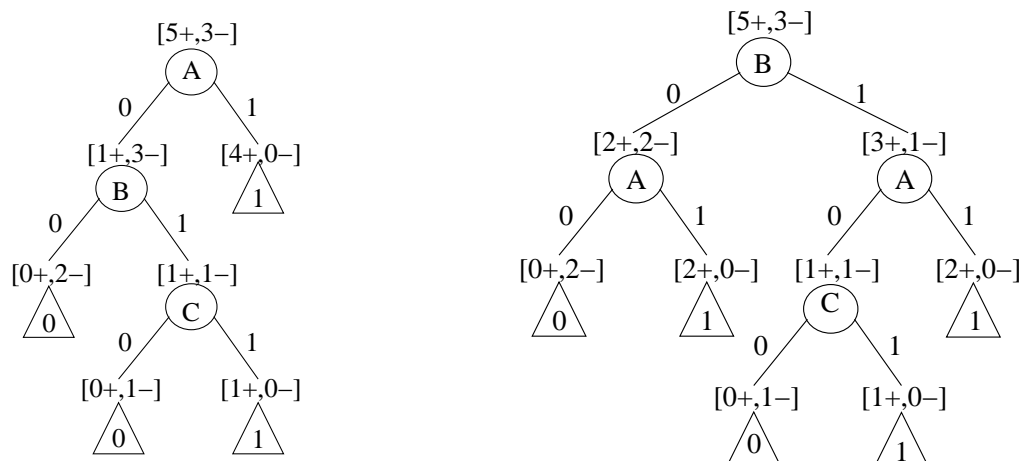
în care simbolul  $H$ , scris uneori însoțit de un argument (așadar, ca  $H(p)$ ), se referă la entropia unei variabile Bernoulli de parametru  $p$ . Mai facem *precizarea* că semnele  $<$  și  $=$  din figura de mai sus se referă de fapt nu [doar] la entropiile condiționale specifice, ci [și] la produsul acestora cu ponderile asociate în mod corespunzător:  $\frac{4}{8}H[1+,3-] = \frac{4}{8}H[3+,1-]$  și respectiv  $\frac{4}{8}H[4+,0-] < \frac{4}{8}H[2+,2-]$ .

**Nodul 1:** Avem de clasificat indivizii care au  $A = 0$  și putem alege între attributele  $B$  și  $C$ . Datorită simetriei, îl putem alege pe oricare dintre ele. Pentru fixare îl alegem pe  $B$ .



**Nodul 2:** La acest punct a mai rămas disponibil doar atributul  $C$ .

Arborele construit de ID3 este cel reprezentat mai jos, în partea stângă:



b. Pentru a vedea dacă arborele construit de algoritmul ID3 este cel optimal, trebuie să reconsiderăm toate deciziile pe care le-am luat în construirea acestuia:

– La nodul 1 al arborelui avem de clasificat exemplele pentru care  $A = 0$ , deci funcția care trebuie reprezentată de subarborele în cauză este  $f' = f[A/0] = 0 \vee (B \wedge C) = B \wedge C$ , funcție care este reprezentată în mod optimal de subarborele construit de ID3. (Notăția  $A/0$  semnifică faptul că variabila logică  $A$  este instanțiată la valoarea 0.) Prin urmare, nu există un arbore mai bun care să reprezinte funcția dată și să aibă în rădăcină atributul  $A$ .

– În rădăcină am ales atributul  $A$  în detrimentul celorlalte două atribute deoarece am demonstrat că aduce cel mai mare câștig de informație. Să vedem ce se întâmplă dacă alegem unul dintre atributele  $B$  sau  $C$ . După cum am discutat mai sus, datorită simetriei, pe oricare dintre cele două l-am alege, arborele rezultat ar avea aceeași formă. Pentru fixare, să-l alegem pe  $B$ .

Subarborele stâng și drept vor trebui să reprezinte funcțiile:

$$f'' = f[B/0] = A \vee (0 \wedge C) = A \vee 0 = A$$

și respectiv

$$f''' = f[B/1] = A \vee (1 \wedge C) = A \vee C$$

Arborele minimal care poate fi construit în aceste circumstanțe este cel reprezentat mai sus în partea dreaptă. Acest arbore are 3 niveluri și 4 noduri, cu un nod în plus față de cel construit de algoritmul ID3.

Putem deci conchide că arborele construit respectând specificațiile algoritmului ID3 este cel optimal.

*Observație:*

Această problemă pune în evidență două modalități de parcurgere a spațiului de versiuni pentru un concept, în particular unul din logica propozițiilor, care este reprezentat cu ajutorul arborilor de decizie. Pe de o parte avem explorarea (incompletă) făcută de algoritmul ID3 care este de tip “greedy”, iar pe de altă parte avem explorarea exhaustivă. Prima strategie de explorare procedează la o căutare „orientată” a soluției (și din această cauză este mai eficientă, dar se va vedea, ca revers, că nu asigură întotdeauna găsirea optimului), iar cea de-a doua strategie de explorare, deși asigură găsirea optimului, nu este utilizabilă în cazurile (frecvente!) în care spațiul de versiuni este foarte mare.

4. (ID3, ca algoritm “greedy”;  
un exemplu când arborele ID3 nu este optimal  
ca număr de noduri și de niveluri)

*prelucrare de Liviu Ciortuz după  
CMU, 2003 fall, T. Mitchell, A. Moore, midterm exam, pr. 9.a*

Fie atributele binare de intrare  $A, B, C$ , atributul de ieșire  $Y$  și următoarele exemple de antrenament:

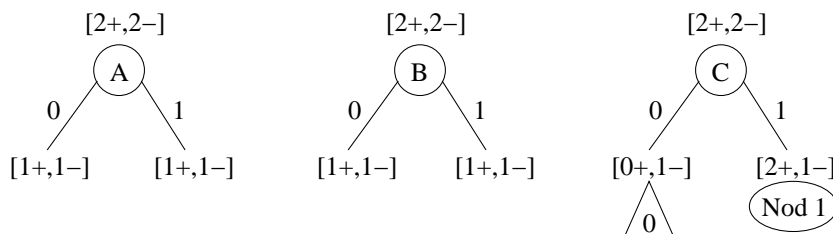
$A$	$B$	$C$	$Y$
1	1	0	0
1	0	1	1
0	1	1	1
0	0	1	0

- a. Determinați arborele de decizie calculat de algoritmul ID3. Este acest arbore de decizie consistent cu datele de antrenament?
- b. Există un arbore de decizie de adâncime mai mică (decât cea a arborelui ID3) consistent cu datele de mai sus? Dacă da, ce concept (logic) reprezintă acest arbore?

**Răspuns:**

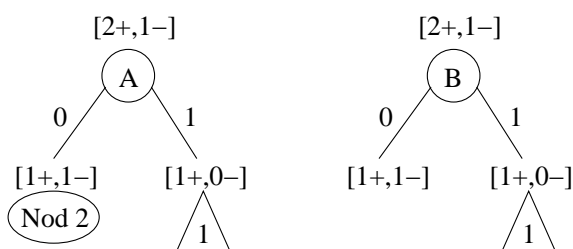
- a. Se construiește arborele de decizie cu algoritmul ID3 astfel:

*Nodul 0* (rădăcina):



Observăm că  $H_{0/A} = H_{0/B} = \frac{2}{4}H[1+, 1-] + \frac{2}{4}H[1+, 1-] = H[1+, 1-] = 1$ , care este valoarea maximă a entropiei [condiționale medii a] unei variabile boolene. Prin urmare,  $H_{0/C}$  nu poate fi decât mai mică sau egală cu  $H_{0/A}$  și  $H_{0/B}$ . Deci vom alege în nodul rădăcină atributul  $C$ .

*Nodul 1:* Avem de clasificat instanțele cu  $C = 1$ , deci alegerea se face între atributele  $A$  și  $B$ .

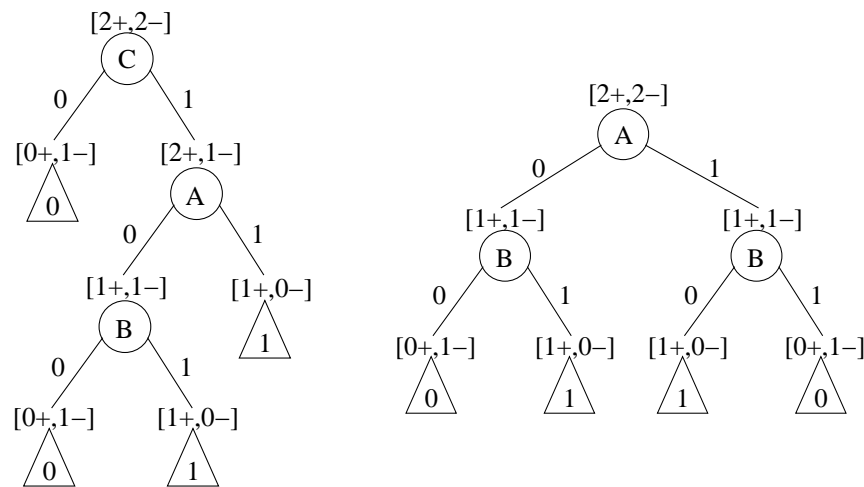


Cele două entropii condiționale medii sunt egale:

$$H_{1/A} = H_{1/B} = \frac{2}{3}H[1+, 1-] + \frac{1}{3}H[1+, 0-]$$

Așadar, putem alege oricare dintre cele două atribute. Pentru fixare, îl alegem pe  $A$ .

*Nodul 2:* La acest nod nu mai avem decât atributul  $B$ , deci îl vom pune pe acesta. Arborele complet este reprezentat în partea stângă:



Prin construcție, algoritmul ID3 este consistent cu datele de antrenament dacă acestea sunt consistente (i.e., necontradictorii). În cazul nostru, se verifică imediat că datele de antrenament sunt consistente.

b. Se observă că atributul de ieșire  $Y$  reprezintă de fapt funcția logică  $A \text{ XOR } B$ . Reprezentând această funcție ca arbore de decizie, vom obține arborele desenat mai sus în partea dreaptă. Acest arbore are cu un nivel mai puțin decât arborele construit cu algoritmul ID3. Prin urmare, arborele ID3 nu este optim din punctul de vedere al numărului de niveluri.

5. (Algoritmul ID3: aplicare pe date inconsistente, “decision stumps”, calculul acurateții)

• ◦ CMU, 2012 fall, T. Mitchell, Z. Bar-Joseph, HW1, pr. 2.ab

Tabelul de mai jos sumarizează situația celor 2201 de pasageri și membri ai echipajului de la bordul vasului Titanic, în urma naufragiului din data de 15 Aprilie 1912. Pentru fiecare combinație de valori ale celor 3 variabile (Clasă, Sex, Vârstă) am indicat în tabel câți oameni au supraviețuit și câți nu au supraviețuit. (*Observație:* Datele originale au patru valori pentru atributul Clasă; am comasat valorile II, III, și Echipaj într-o singură valoare, denumită „Inferioară“.)

Clasa	Sexul	Vârsta	Supraviețuitori		
			Nu	Da	Total
I	Masculin	Copil	0	5	5
I	Masculin	Adult	118	57	175
I	Feminin	Copil	0	1	1
I	Feminin	Adult	4	140	144
Inferioară	Masculin	Copil	35	24	59
Inferioară	Masculin	Adult	1211	281	1492
Inferioară	Feminin	Copil	17	27	44
Inferioară	Feminin	Adult	105	176	281
Total			1490	711	2201

Pentru a vă ușura calculele pe care va trebui să le faceți, am făcut noi totalurile pentru fiecare variabilă:

Clasa	Supraviețuitori		
	Nu	Da	Total
I	122	203	325
Inferioară	1368	508	1876

Sexul	Supraviețuitori		
	Nu	Da	Total
Masculin	1364	367	1731
Feminin	126	344	470

Vârsta	Supraviețuitori		
	Nu	Da	Total
Copil	52	57	109
Adult	1438	654	2092

a. Folosind un arbore de decizie, dorim să prezicem variabila de ieșire  $Y$  (Supraviețuitor), pornind de la atributele de intrare  $C$  (Clasa),  $S$  (Sexul),  $V$  (Vârsta). Utilizați criteriul câștigului de informație pentru a alege care dintre aceste trei atribute  $C$ ,  $S$  sau  $V$  trebuie să fie folosit în nodul-rădăcină al arborelui de decizie.

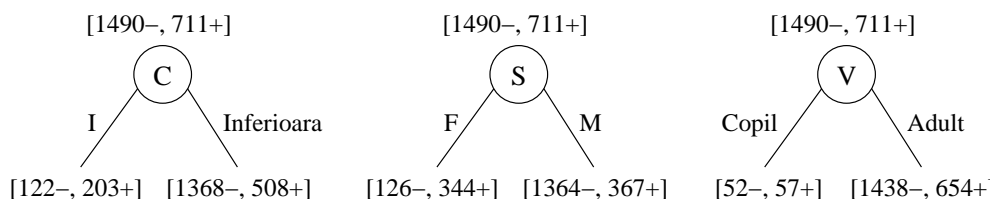
De fapt, ce vi se cere este să învățați un arbore de decizie de adâncime 1 care folosește doar atributul din rădăcină pentru a clasifica datele. (Astfel de arbori de decizie de adâncime 1 sunt adesea numiți în terminologia de limbă engleză “decision stumps”.) Parcurgeți toate etapele rezolvării, redând inclusiv calculele pentru câștigul de informație al fiecărui atribut.

b. Care este acuratețea [medie] obținută pe datele de antrenament de către arborele de decizie cu adâncime 1 de la punctul precedent?

c. Dacă ați crea un arbore de decizie care folosește toate cele trei variabile, care ar fi acuratețea lui [medie] pe datele de antrenament? (*Observație:* Nu trebuie neapărat să creați arborele de decizie pentru a afla răspunsul!)

#### Răspuns:

a. Totalurile care au fost furnizate în enunț pentru fiecare dintre variabilele  $C$ ,  $S$  și  $V$  ne servesc foarte bine pentru a crea rapid cei trei “decision stumps”:



Analizând datele conform figurii de mai sus, se poate „intui” că atributul  $S$  va avea un câștig de informație (în raport cu atributul de ieșire  $Y$  – *Supraviețuitor*) mai bun decât al celorlalte două atribute de intrare ( $C$  și  $V$ ). Intuiția se verifică făcând calculele:

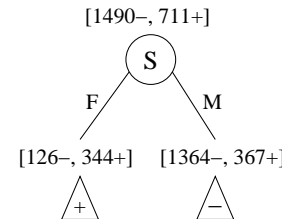
$$\begin{aligned}
 IG(Y, C) &= H[1490-, 711+] - \left( \frac{325}{2201} H[122-, 203+] + \frac{1876}{2201} H[1368-, 508+] \right) \\
 &= 0.048501 \\
 IG(Y, S) &= H[1490-, 711+] - \left( \frac{470}{2201} H[126-, 344+] + \frac{1731}{2201} H[1364-, 367+] \right) \\
 &= 0.142391
 \end{aligned}$$

$$\begin{aligned}
 IG(Y, V) &= H[1490-, 711+] - \left( \frac{109}{2201} H[52-, 57+] + \frac{2092}{2201} H[1438-, 654+] \right) \\
 &= 0.006411.
 \end{aligned}$$

Deci, într-adevăr, câștigul maxim de informație se obține pentru atributul  $S$ .

b. Arborele de decizie de adâncime 1 care are în nodul rădăcină atributul  $S$  este cel din figura alăturată. Acuratețea [medie  $a$ ] acestui arbore de decizie este:

$$\frac{470}{2201} \cdot \frac{344}{470} + \frac{1731}{2201} \cdot \frac{1364}{1731} = \frac{344 + 1364}{2201} = \frac{1708}{2201} = 0.776.$$



c. Se poate constata imediat că arborele ID3 produs pe datele din această problemă va avea 8 noduri-frunză, iar în fiecare dintre aceste noduri-frunză se va asigura câte una dintre mulțimile descrise (pe linie) în coloanele 4 și 5 ale tabelului principal din enunț:  $[0-, 5+]$ ,  $[118-, 57+]$ ,  $\dots$ ,  $[17-, 27+]$ ,  $[105-, 176+]$ . Decizia care va fi luată în fiecare nod-frunză este dictată de votul majoritar, și anume:  $+$ ,  $-$ ,  $\dots$ ,  $+$  și respectiv  $+$ .

Putem calcula acuratețea [medie] astfel:

$$\frac{5 + 118 + 1 + 140 + 35 + 1211 + 27 + 176}{2201} = \frac{1713}{2201} = 0.778.$$

Se observă că se produce (din păcate) o creștere foarte mică în raport cu acuratețea celui mai bun "decision stump": doar 0.002.

6. (Algoritmul ID3: cazul când există repetiții și inconsistențe în datele de antrenament; o margine superioară pentru eroarea la antrenare în funcție de numărul de valori ale variabilei de ieșire)  
*CMU, 2002 fall, Andrew Moore, midterm exam, pr. 1.fg*

Presupunem că învățăm un arbore de decizie care să prezică atributul de ieșire  $Z$  pornind de la atributele de intrare  $A, B, C$ . Se folosesc datele de antrenament din tabelul alăturat.

a. Care va fi eroarea la antrenare pe acest set de date? Exprimați răspunsul sub forma fracției de înregistrări care vor fi clasificate eronat ( $n/12$ ).

b. Considerăm un arbore de decizie construit pe un set arbitrar de date. Dacă atributul de ieșire este cu valori discrete și poate lua  $k$  valori distincte, care este eroarea de antrenare maximă (exprimată ca fracție)?

$A$	$B$	$C$	$Z$
0	0	0	0
0	0	1	0
0	0	1	0
0	1	0	0
0	1	1	0
0	1	1	1
1	0	0	0
1	0	1	1
1	1	0	1
1	1	0	1
1	1	1	0
1	1	1	1

Răspuns:

a. Este ușor de observat că datele de antrenament conțin „inconsistențe” (contradicții, relativ la etichetare), și anume la exemplele  $(0, 1, 1)$  și  $(1, 1, 1)$ .

Fiecare dintre aceste exemple sunt etichetate o dată cu 0 și altă dată cu 1. Prin urmare, jumătate din aceste exemple vor fi clasificate eronat de arborele învățat de către algoritmul ID3. Eroarea la antrenare va fi deci  $\frac{2}{12}$ .

b. Vom analiza pe rând mai multe cazuri, care sunt din ce în ce mai generale.

Cazul *i*: Mai întâi vom calcula eroarea la antrenare pentru cazul în care setul de date de antrenament este compus din  $k$  instanțe care sunt identice ca tupluri de valori pentru atributele ce le caracterizează, dar sunt clasificate pe rând cu fiecare din cele  $k$  valori posibile ale atributului de ieșire. Este evident că arborele de decizie învățat va clasifica eronat  $k - 1$  instanțe. Așadar, în acest caz, eroarea la antrenare va fi

$$E = \frac{k-1}{k}$$

Cazul *ii*: Aceeași eroare [la antrenare] ca mai sus se va înregistra dacă în locul fiecărei instanțe dintre cele considerate la cazul precedent vom avea  $l$  instanțe identice, inclusiv în ce privește eticheta. (În total sunt  $kl$  instanțe de antrenament.)

$$E = \frac{(k-1) \cdot l}{k \cdot l} = \frac{k-1}{k}$$

Cazul *iii*: Dacă relaxăm condiția de mai sus considerând  $l_1, l_2, \dots, l_k$  instanțe identice, iar  $l = \max_{i=1}^k l_i$ , este imediat că eroarea maximă se va atinge în cazul  $l_1 = l_2 = \dots = l_k = l$ , și va avea aceeași valoare ca mai sus.

Așadar, în continuare vom putea renunța la a considera factorul de multiplicare  $l$ , fără ca prin aceasta să restrângem generalitatea raționamentului.

Cazul *iv*: Fie  $n$  exemple de antrenament (instanțe etichetate) dintre care  $d$  sunt distincte (ca tupluri de valori ale atributelor de intrare). Fie  $k_1, k_2, \dots, k_d$  numărul de instanțe etichetate pentru fiecare caz distinct în parte din cele  $d$ . Atunci vom avea:

$$k_1 \leq k, k_2 \leq k, \dots, k_d \leq k \Rightarrow n = k_1 + k_2 + \dots + k_d \leq k \cdot d \text{ deci } n \leq k \cdot d$$

Eroarea maximă la antrenare va fi dată de formula

$$E = \frac{(k_1 - 1) + (k_2 - 1) + \dots + (k_d - 1)}{n} = \frac{n - d}{n}$$

Avem:

$$E \leq \frac{k-1}{k} \Leftrightarrow \frac{n-d}{n} \leq \frac{k-1}{k} \Leftrightarrow k \cdot n - k \cdot d \leq k \cdot n - n \Leftrightarrow k \cdot d \geq n \text{ (adev.)}$$

Prin urmare, eroarea maximă la antrenare care poate fi atinsă atunci când atributul de ieșire poate lua  $k$  valori distincte este  $\frac{k-1}{k}$ .

7. (ID3, aspecte computaționale:  
influența atributelor duplicate, respectiv a  
instanțelor de antrenament duplicate  
asupra arborelui ID3 rezultat)

*CMU, 2009 spring, Ziv Bar-Joseph, final exam, pr. 3.1-2*

Se dorește construirea unui arbore de decizie pentru  $n$  vectori, cu  $m$  atribute.

a. Să presupunem că există  $i$  și  $j$  astfel încât pentru TOȚI vectorii  $X$  din datele de antrenament, aceste atribute au valori egale ( $x_i = x_j$  pentru toți vectorii, unde  $x_i$  este valoarea atributului  $i$  în vectorul  $X$ ). Să presupunem de asemenea că în cazul în care ambele atribute duc la același câștig de informație vom folosi atributul  $i$ . Ștergerea atributului  $j$  din datele de antrenament poate schimba arborele de decizie obținut? Explicați pe scurt.

b. Să presupunem că există în mulțimea de antrenament doi vectori egali  $X$  și  $Z$  (adică, toate atributele lui  $X$  și  $Z$  sunt exact la fel, inclusiv etichetele). Ștergerea vectorului  $Z$  din datele de antrenament poate schimba arborele de decizie obținut? Explicați pe scurt.

Răspuns:

a. Nu, îndepărtarea atributului  $j$  nu schimbă arborele de decizie, deoarece atributele  $i$  și  $j$  conduc la valori egale pentru câștigul de informație în fiecare nod al arborelui.

b. Da, în acest caz arborele de decizie se poate schimba, fiindcă entropia condițională — care se calculează în fiecare nod pentru a determina atributul cu câștigul de informație cel mai mare — depinde de numărul de instanțe de antrenament luate în considerare.

8. (Arbori de decizie: o margine superioară pentru  
numărul de noduri frunză, în funcție de  
numărul atributelor și numărul de exemple)

*CMU, 2005 fall, T. Mitchell, A. Moore, midterm exam, pr. 2.d*

Presupunem că învățăm un arbore de decizie folosind un set de  $R$  instanțe de antrenament descrise de  $M$  atribute de intrare având valori binare.

Care este numărul maxim posibil de noduri frunză din arborele de decizie, presupunând că fiecărui nod frunză îi este asociat măcar un exemplu de antrenament? Încercuiți unul din răspunsurile de mai jos; justificați alegerea făcută.

$R, \log_2(R), R^2, 2^R, M, \log_2(M), M^2, 2^M,$   
 $\min(R, M), \min(R, \log_2(M)), \min(R, M^2), \min(R, 2^M),$   
 $\min(\log_2(R), M), \min(\log_2(R), \log_2(M)), \min(\log_2(R), M^2), \min(\log_2(R), 2^M),$   
 $\min(R^2, M), \min(R^2, \log_2(M)), \min(R^2, M^2), \min(R^2, 2^M),$   
 $\min(2^R, M), \min(2^R, \log_2(M)), \min(2^R, M^2), \min(2^R, 2^M),$   
 $\max(R, M), \max(R, \log_2(M)), \max(R, M^2), \max(R, 2^M),$   
 $\max(\log_2(R), M), \max(\log_2(R), \log_2(M)), \max(\log_2(R), M^2), \max(\log_2(R), 2^M),$   
 $\max(R^2, M), \max(R^2, \log_2(M)), \max(R^2, M^2), \max(R^2, 2^M),$   
 $\max(2^R, M), \max(2^R, \log_2(M)), \max(2^R, M^2), \max(2^R, 2^M).$



**Răspuns:**

Notăm cu  $\max_{frunze}$  valoarea căutată. Trebuie luate în considerare două aspecte:

- (a) fiecare nod frunză trebuie să clasifice măcar un exemplu de antrenament  
 $\Rightarrow$  nu putem să avem mai multe frunze decât exemple de antrenament  
 $\Rightarrow \max_{frunze} \leq R$
- (b) fiecare atribut poate fi testat o singură dată pe un drum de la rădăcină la o frunză oarecare  $\Rightarrow$  arborele obținut va avea adâncimea cel mult  $M$   
 $\Rightarrow \max_{frunze} \leq 2^M$

$$\left. \begin{array}{l} \max_{frunze} \leq R \\ \max_{frunze} \leq 2^M \end{array} \right\} \Rightarrow \max_{frunze} \leq \min(R, 2^M)$$

Această valoare poate fi atinsă: putem să luăm, de exemplu, cazul trivial când avem o singură instanță de antrenament. Prin urmare, avem  $\max_{frunze} = \min(R, 2^M)$ .

## 3.2 Arbori de decizie — Probleme propuse

28. (Arbori de decizie; optimalitate, ca număr minim de noduri)

\*

Reprezentați arborele/arborii de decizie care are/au numărul minim posibil de noduri (de test) și corespunde/corespund funcției boolene  $(A \text{ XOR } B) \wedge C$  definită peste atributele boolene  $A, B$  și  $C$ .

29. (Expresivitatea arborilor de decizie: un rezultat privind funcțiile boolene)

Orice funcție booleană (care primește  $n$  argumente din mulțimea  $\{0, 1\}$  și întoarce un element din mulțimea  $\{0, 1\}$ ) poate fi reprezentată cu ajutorul unui arbore de decizie. Adevărat sau fals?

În cazul afirmativ, explicați succint cum anume poate fi construit arborele de decizie respectiv.

În cazul negativ, dați un exemplu de funcție booleană pentru care nu se poate construi un arbore de decizie consistent cu funcția respectivă.

30. (Calcularea câștigului de informație pe “decision stumps”)

□ • CMU, 2013 fall, W. Cohen, E. Xing, *Sample Questions*, pr. 4

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 colectat în total nouă *instance*/exemple, descrise cu ajutorul a două *trăsături* (văzute în cele ce urmează 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):

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

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:

$N$	3	5	7
$\log_2 N$	1.5850	2.3219	2.8073

31. (Implementare: entropie, entropie condițională specifică, entropie condițională medie, câștig de informație)

□ *Liviu Ciortuz, 2016*

Folosind limbajul de programare pe care-l preferați, implementați un program care, pornind de la o structură de date de tip compas de decizie (engl., decision stump), calculează entropia, entropiile condiționale specifice, entropia condițională medie, precum și câștigul de informație aferent.

În mod concret, programul va primi ca *input*

- $m$  — numărul de valori posibile ale etichetei / atributului de ieșire (în mod implicit, se va considera  $m = 2$ );
- $n$  — numărul de valori ale atributului (notat mai jos cu  $A$ ) în raport cu care se face partiționarea mulțimii de instanțe asociate nodului-rădăcină al compasului de decizie (valoarea implicită:  $n = 2$ );
- partițiile (de fapt, count-urile) corespunzătoare nodurilor descendente. Pornind de la aceste partiții, programul va calcula partiția asociată nodul-rădăcină al compasului de decizie.

De *exemplu*, pentru primul compas de decizie de la problema 30, inputul va avea forma  $[3, 2]$ ,  $[2, 2]$ , în vreme ce pentru al doilea compas de decizie va fi  $[5, 1]$ ,  $[0, 3]$ .

Programul va calcula și apoi va afișa

- entropia atributului / variabilei de ieșire (notată aici cu  $Y$ );
- entropiile condiționale specifice pentru fiecare descendent din nodul-rădăcină;
- entropia condițională medie a atributului  $A$ ;
- câștigul de informație al atributului [de ieșire]  $Y$  în raport cu atributul [de intrare]  $A$ .

32. (Algoritmul ID3: aplicare pe expresii booleene; exploatarea simetriilor operațiilor  $\vee, \wedge$  în alegerea nodurilor; analiza „optimalității” arborelui ID3)

\* *prelucrare de Liviu Ciortuz după Tom Mitchell, “Machine Learning”, 1997, ex. 3.1.d*

Considerăm următoarea funcție booleană:  $(A \wedge B) \vee (C \wedge D)$ . Valorile pe care le ia această funcție, calculate conform diferitelor valori de adevăr atribuite variabilelor/atributelor  $A, B, C$  și  $D$  sunt cele cunoscute din logica propozițiilor. Dorim însă să reprezentăm această funcție ca arbore de decizie.

a. Aplicați algoritmul ID3 acestei funcții.

*Observație:* Dacă exploatați simetriile, este nevoie doar de puține calcule, altfel vă complicați în mod inutil.

b. Arborele ID3 obținut la punctul precedent este optimal?

Alfel spus, puteți găsi alt arbore de decizie de adâncime mai mică sau cu număr mai mic de noduri (de test) pentru această funcție? (Țineți cont că în fiecare nod al unui arbore de decizie se poate testa un singur atribut.)

33.

(Algoritmul ID3: aplicare  
analiza „optimalității“ arborelui ID3)

• \* CMU, 2005 spring, C. Guestrin, T. Mitchell, midterm, pr. 4

Agencia spațială NASA dorește să distingă între marțieni (M) și pământeni (H) folosind următoarele caracteristici:  $Green \in \{N, Y\}$ ,  $Legs \in \{2, 3\}$ ,  $Height \in \{S, T\}$ ,  $Smelly \in \{N, Y\}$ .

Datele de antrenament de care dispunem sunt prezentate în tabelul alăturat.

	Species	Green	Legs	Height	Smelly
1	M	N	3	S	Y
2	M	Y	2	T	N
3	M	Y	3	T	N
4	M	N	2	S	Y
5	M	Y	3	T	N
6	H	N	2	T	Y
7	H	N	2	S	N
8	H	N	2	T	N
9	H	Y	2	S	N
10	H	N	2	T	Y

a. Învățați un arbore de decizie folosind algoritmul ID3 și trasați arborele respectiv.

b. Descrieți conceptul M (marțian) ca un set de reguli conjunctive din logica propozițiilor. Spre exemplu:

```

if  $Green = Y$  and  $Legs = 2$  and  $Height = T$  and  $Smelly = N$  then M;
else
if ... then M; else H.

```

c. Soluția de la punctul b de mai sus folosește cel mult 4 atribute în fiecare conjuncție. Găsiți un set de reguli conjunctive care folosesc doar 2 atribute pentru fiecare conjuncție, păstrând însă eroarea la antrenare zero. Această ipoteză mai simplă poate fi reprezentată ca un arbore de decizie de adâncime 2? Justificați răspunsul.

34.

(Algoritmul ID3: aplicare;  
cazul instanțelor de antrenament cu multiple apariții)

CMU, 2010 fall, Aarti Singh, HW2, pr. 5.1

Tabelul alăturat descrie instanțe (înregistrări) pozitive și instanțe negative pentru persoane cărora banca le-a acordat (sau nu le-a acordat) un card de credit. Fiecare linie din tabel indică niște combinații de valori observate pentru atribuțiile considerate ( $Gender$ ,  $Income$  și  $Approved$ ) și de câte ori a fost înregistrată respectiva combinație de valori.

$Gender$	$Income$	$Approved$	$Counts$
F	Low	+	10
F	High	+	95
M	Low	+	5
M	High	+	90
F	Low	-	80
F	High	-	20
M	Low	-	120
M	High	-	30

De exemplu, (F, Low, +) a apărut de 10 ori, iar (F, Low, -) de 80 de ori.

- Calculați entropia atributului *Approved* pe acest set de date de antrenament (folosind logaritmul cu baza 2).
- Calculați de asemenea câștigurile de informație  $IG(Approved, Gender)$  și  $IG(Approved, Income)$ .
- Desenați un arbore de decizie produs de către algoritmul ID3 (fără post-pruning) pe baza acestui set de date de antrenament.

35. (“Decision stump” produs de ID3: raționament calitativ pe un exemplu simplu)

• ◦ CMU, 2010 fall, Ziv Bar-Joseph, midterm exam, pr. 5.a

Vrem să construim un arbore de decizie care să ne ajute să prezicem întârzierile avioanelor. Timp de câteva luni am colectat informații, iar un rezumat al acestora este prezentat în tabelul următor:

Atribut	Valoare = <i>Da</i>		Valoare = <i>Nu</i>	
	#Zboruri		#Zboruri	
	amânate	neamânate	amânate	neamânate
Ploaie	30	10	10	30
Vânt	25	15	15	25
Vara	5	35	35	5
Iarna	20	10	20	30
Ziua	20	20	20	20
Noaptea	15	10	25	30

- Pe baza acestui tabel precizați ce atribut ar trebui să fie pus în rădăcina arborelui de decizie. (Justificați pe scurt, nu este necesar să calculați valorile exacte pentru câștigul de informație.)
- Pe baza aceluiași tabel, precizați care dintre atribute ar trebui să apară pe al doilea nivel (nivelul de sub rădăcină) al arborelui de decizie. (Justificați pe scurt, nu este necesar să calculați valorile exacte pentru câștigul de informație.)

36. (Clasificare ternară: “decision stump” produs de ID3, pe date care conțin duplicări și „zgomote”)

• \*

Presupunem că se dau șase date de antrenament (precizate în tabel) pentru o problemă de clasificare cu două atribute binare și trei clase  $Y \in \{1, 2, 3\}$ . Se va crea un arbore ID3, bazat pe câștigul de informație.

- Calculați câștigul de informație atât pentru  $X_1$  cât și pentru  $X_2$ . Se va folosi aproximarea  $\log_2 3 = 19/12$  și se va scrie câștigul de informație sub formă de fracții.

$X_1$	$X_2$	$Y$
1	1	1
1	1	1
1	1	2
1	0	3
0	0	2
0	0	3

b. Pe baza rezultatelor anterioare, ce atribut va fi folosit pentru primul nod al arborelui ID3? Desenați arborele de decizie care rezultă folosind doar acest singur nod. Etichetați corespunzător nodul, ramurile și eticheta prevăzută în fiecare frunză.

c. Cum va clasifica acest arbore instanța determinată de  $X_1 = 0$  și  $X_2 = 1$ ?

37. (O aproximare a numărului de instanțe greșit clasificate care au fost asignate la un nod frunză dintr-un arbore ID3)

• ◦ *CMU, 2003 fall, Andrew Moore, midterm exam, pr. 9.b*

Învățăm un arbore de decizie folosind un set de date de antrenament cu atributul de ieșire (*class*) având valorile 0 sau 1.

Presupunând că pentru un nod frunză  $l$  din acest arbore,

- există  $M$  instanțe de antrenament asignate la acel nod, iar
- entropia sa este  $H$ ,

schițați un algoritm simplu care ia ca valori de intrare  $M$  și  $H$  și furnizează la ieșire numărul de exemple de antrenament clasificate greșit de către nodul frunză  $l$ .

*Sugestie:* Folosiți o aproximare simplă (polinomială) pentru funcția entropie  $H(p)$ .

38. (Algoritmul ID3: eroarea la antrenare)

• ◦ *CMU, 2003 fall, Andrew Moore, HW1, pr. 2.1*

Un student mi-a spus următoarele:

- el poate să construiască un set de instanțe cu atributele de intrare discrete și atributul de ieșire binar;
- mie îmi dă voie să aleg o parte din acest set de instanțe (dar nu toate!) pentru a antrena un arbore de decizie;
- indiferent de modul cum mi-aș alege datele de antrenament din setul construit de el, eroarea de clasificare pe care arborele de decizie (obținut în urma antrenării) o va face pe instanțele care nu au fost incluse în setul de antrenament va fi de cel puțin 50%.

Credeți că studentul are dreptate? Explicați de ce sau dați un exemplu.