

CONSTRUCȚIA ARBORILOR DE CLASIFICARE - ALGORITMUL ID3

ing. Lorina Negreanu

Universitatea "Politehnica" București

Rezumat: Achiziția automată de cunoștințe este unul dintre principalele domenii de aplicabilitate a conceptelor și tehnicilor de învățare automată. O metoda eficientă și deja foarte populară pentru inducerea regulilor de clasificare din exemple este algoritmul ID3 (Quinlan). Prezentul articol descrie algoritmul ID3 și prezintă o extindere a utilizării sale pentru probleme cu caracter nedeterminist. De asemenea sunt prezentate pe scurt câteva optimizări ale algoritmului, respectiv algoritmi ID3 Generalizat, ID4 și ID5.

Cuvinte cheie: arbore de clasificare, atribut, clasă, exemple, mesaj, informație, grad de certitudine, nedeterminism, bază de cunoștințe, diagnoză.

1. Prezentarea algoritmului ID3

1.1. Principiul generării inductive a arborilor de clasificare

Algoritmul ID3, un descendent al algoritmului CLS [2], pornește de la o mulțime de obiecte descrise printr-o mulțime fixă de proprietăți și construiește un arbore de decizie (clasificare) care poate clasifica toate obiectele din mulțime.

Un obiect este caracterizat printr-o mulțime fixă de atribute (proprietăți), fiecare atribut având la rândul său o mulțime de valori posibile. Descrierea unui obiect (exemplu) constă într-o succesiune de valori de atribut. Dacă C este o mulțime de obiecte, o regulă de clasificare sub forma unui arbore de decizie poate fi construită astfel:

- dacă C este mulțimea vidă ea este asociată arbitrar cu orice clasă;
- dacă toate obiectele din C aparțin aceleiași clase, atunci arborele de decizie constă doar dintr-o

frunză ce are numele clasei respective;

-altfel, se selectează un atribut și se împarte C în submulțimi disjuncte C_1, C_2, \dots, C_n , unde C_i conține acei membri din C care au valoarea i a atributului selectat; pentru fiecare dintre aceste submulțimi se aplică strategia prezentată.

În final, se generează un arbore în care frunzele sunt etichetate cu nume de clase, iar nodurile interioare specifică atributul care trebuie testat, având arce care pleacă din nod pentru fiecare valoare posibilă a atributului respectiv.

Pentru a ilustra cele prezentate, se consideră atributele *înălțime* având valorile {înalt, scund}, *culoare_păr* având valorile {șaten, blond, roșcat}, *culoare_ochi* având valorile {verde, negru}. Clasele posibile sunt {+, -}. Mulțimea exemplilor (C) este următoarea:

$C =$

<i>înălțime</i>	<i>culoare_păr</i>	<i>culoare_ochi</i>	<i>clasa</i>
scund	blond	verde	+
înalt	blond	negru	-
înalt	roșcat	verde	+
scund	șaten	verde	-
înalt	șaten	verde	-
înalt	blond	verde	+
înalt	șaten	negru	-
scund	blond	negru	-

Pentru a construi arborele de decizie trebuie selectat un atribut care să fie rădăcina arborelui. Fie acesta, atributul *culoare_păr*. Arborele care rezultă pentru acest pas este prezentat în figura 1.

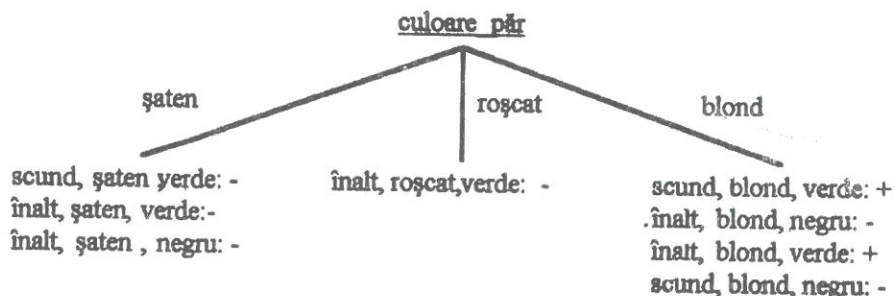


Figura 1. Arbore de decizie parțial

Submulțimile corespunzătoare valorilor *șaten* și *roșcat* conțin obiecte care aparțin aceleiași clase, deci ele nu mai trebuie prelucrate. Pentru ramura corespunzătoare valorii de atribut *blond*, trebuie

ales un alt atribut pentru a fi testat. Fie acesta *culoare_ochi*. Arborele care se obține este prezentat în figura 2.

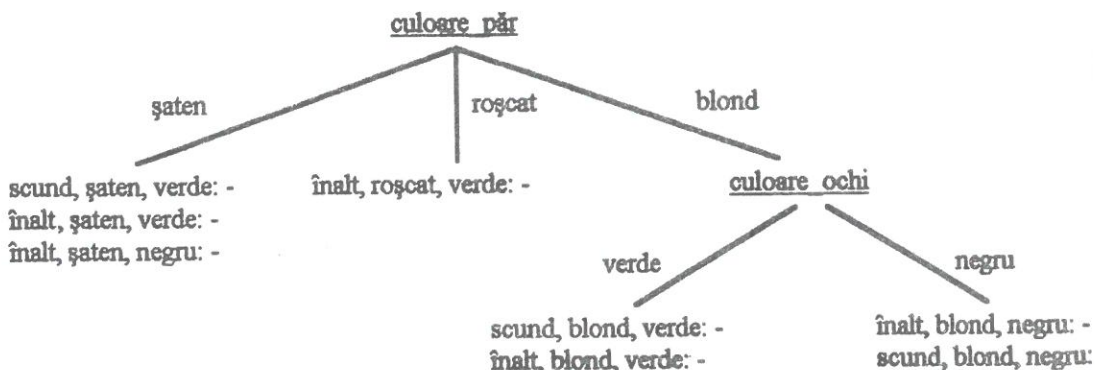


Figura 2. Arbore de clasificare complet

Toate submulțimile conțin obiecte care aparțin aceleiași clase, deci pot fi înlocuite cu noduri

etichetate cu numele claselor, obținându-se arborele prezentat în figura 3.

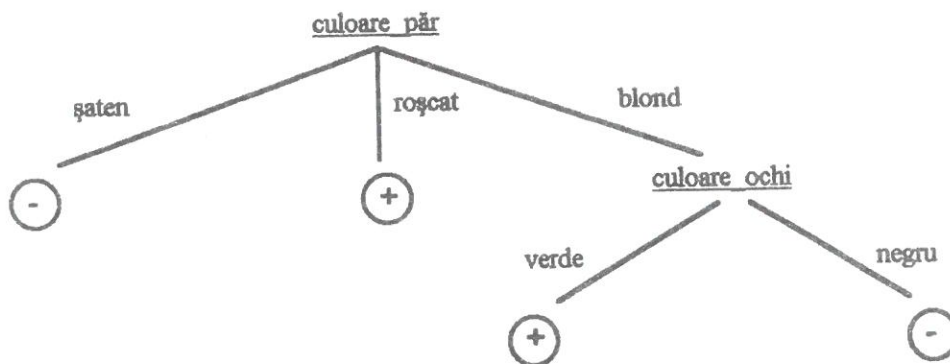


Figura 3. Arbore de clasificare

Pentru a clasifica un obiect, se pornește din rădăcina arborelui de decizie și se urmează calea dictată de valorile proprii obiectului respectiv, până se ajunge la o frunză. Așa cum se observă, clasificarea unui obiect implică evaluarea unui număr mic de atribute, respectiv a celor care se află pe calea de la rădăcină la frunza corespunzătoare. Arborele construit trebuie să poată clasifica și obiecte care nu au fost în mulțimea inițială de exemple. Obiectele care au aceleași valori de atribut, dar aparțin unor clase diferite, reprezintă inconsistențe ale mulțimii de exemple și desigur ele nu pot fi clasificate.

În general, dimensiunea arborelui de clasificare depinde de ordinea în care sunt alese atributele pentru a fi evaluate. Deci, problema care se pune este următoarea: având o mulțime de exemple și un număr de arbori de clasificare diferiți, care este cel care are cele mai multe șanse să clasifice corect obiecte care nu aparțin mulțimii inițiale?

Algoritmul ID3 pornește de la premisa că arborele optim este cel mai "simplu" arbore care satisface mulțimea de exemple. Rațiunea acestei

euristici este principiul enunțat de logicianul medieval William din Occam în 1324 (*Occam's Razor*): este inutil să faci cu mai mult, ceea ce poți să faci cu mai puțin.

Dacă presupunem că mulțimea inițială de exemple este suficientă pentru a construi o generalizare validă, atunci problema se reduce la diferențierea proprietăților necesare.

1.2. Evaluarea unei proprietăți

O proprietate poate fi evaluată prin prisma contribuției informaționale pe care o aduce pentru procesul de clasificare. De exemplu, dacă dorim să determinăm specia unui animal, proprietatea "are pene", contribuie cu o anumită cantitate de informație la atingerea scopului. ID3 măsoară câștigul potențial de informație al fiecărei proprietăți, punând-o ca rădăcină a subarborelui curent și alege proprietatea pentru care câștigul este maxim. Baza matematică pentru această evaluare este teoria informațională a lui Shannon(1984).

Un mesaj poate fi considerat ca o instanță într-un univers de mesaje posibile. Actul transmiterii unui mesaj este similar selectării unui mesaj din univers. Din acest punct de vedere conținutul de informație al unui mesaj poate fi considerat dependent de mărimea universului și de frecvența de apariție a mesajului. Shannon a formalizat aceste intuiții definind cantitatea de informație dintr-un mesaj ca o funcție a probabilității de apariție a fiecărui mesaj.

Fie $M = \{m_1, m_2, \dots, m_n\}$ un univers de mesaje și $p(m_i)$, $i = 1, n$ probabilitatea de apariție a fiecărui mesaj. Cantitatea de informație a unui mesaj este dată de relația:

$$I(M) = \sum_{i=1}^n -p(m_i) \log_2(p(m_i))$$

Algoritmul ID3 utilizează teoria informației pentru a selecta atributul (testul) care aduce cel mai mare câștig de informație pentru clasificarea mulțimii de exemple. Pentru un anumit test, câștigul de informație obținut, dacă testul respectiv este făcut în rădăcina arborelui, este egal cu diferența dintre informația din arbore și informația necesară pentru realizarea clasificării după ce testul a fost efectuat. Cantitatea de informație necesară pentru a construi arborele este definită ca o medie a cantității de informație din fiecare subarbore ponderată cu procentul exemplelor prezente în subarbore.

Fie C o mulțime de exemple. Dacă proprietatea P , care are n valori, devine rădăcina subarborelui curent, aceasta va partiționa mulțimea C în submulțimile $\{C_1, C_2, \dots, C_n\}$ corespunzătoare valorilor atributului. Cantitatea așteptată de informație, necesară construirii arborelui a cărui rădăcină este P , este:

$$E(P) = \sum_{i=1}^n \frac{|C_i|}{|C|} I(C_i)$$

Câștigul obținut prin testarea proprietății P este diferența dintre conținutul de informație al arborelui și cantitatea de informație necesară pentru construcția arborelui:

$$c(P) = I(C) - E(P)$$

Ca rădăcină a arborelui curent este selectată proprietatea care aduce cel mai mare câștig de informație.

În cazul exemplului prezentat,

$I(C) = -3/8 \log_2 3/8 - 5/8 \log_2 5/8 = 0,954$
Pentru a decide care atribut să devină rădăcina arborelui, trebuie calculată cantitatea de informație estimată pentru fiecare atribut în parte.

În cazul atributului *înălțime*, pentru valoarea *înalt*:

$$I(C_1) = -2/5 \log_2 2/5 - 3/5 \log_2 3/5 = 0.971$$

iar pentru valoarea *scund*:

$$I(C_2) = -1/3 \log_2 1/3 - 2/3 \log_2 2/3 = 0.918$$

Cantitatea estimată de informație, este:

$$E(\text{inaltime}) = 5/8 * 0.971 + 3/8 * 0.918 = 0.951$$

iar câștigul obținut prin testarea acestui atribut este:

$$0.954 - 0.951 = 0.003,$$

practic neglijabil.

În cazul celui de al doilea atribut, *culoare_păr*, valorile sunt *roșcat*, *blond*, *șaten*, iar cantitatea de informație estimată, calculată analog, este:

$$E(\text{culoare_păr}) = 0.5$$

Câștigul obținut prin testarea acestui atribut este:

$$0.954 - 0.5 = 0.454$$

Printr-un calcul analog, câștigul de informație în cazul testării atributului *culoarea_ochi* este 0.347. Principiul maximi-zării câștigului de informație va determina algoritmul ID3 să selecteze atributul *culoare_păr* ca rădăcină a arborelui de decizie.

Metoda descrisă pentru construcția arborilor de clasificare presupune că operațiile asupra mulțimii C de obiecte (cum ar fi determinarea mulțimii corespunzătoare valorii de atribut A_i a atributului A) pot fi făcute eficient, ceea ce practic înseamnă că mulțimea C trebuie păstrată în memoria internă. Ce se întâmplă însă dacă dimensiunea mulțimii C nu permite acest lucru? Soluția adoptată de ID3 se bazează pe o rafinare incrementală a arborilor de clasificare.

Se selectează aleator o submulțime de instanțe, numită *ferastră*. Pentru această submulțime se construiește arborele de decizie. Dintre exemplele care au rămas în afara ferestrei se selectează cele care nu pot fi clasificate cu arborele construit și se formează o nouă fereastră care include și aceste exemple. Se construiește un nou arbore de decizie. Procesul se reia până când se construiește un arbore pentru care nu mai există excepții, deci care clasifică, în mod corect, toate exemplele din mulțimea inițială C .

Pentru formarea unei noi ferestre, în [4] sunt prezentate două metode. Prima metodă mărește fereastra curentă prin adăugarea unui număr de excepții până la o limită dată. Cea de a doua metodă încearcă să identifice obiectele "cheie" din fereastra curentă și să le înlocuiască pe celelalte cu excepțiile determinate, păstrând astfel dimensiunea ferestrei constantă. Ambele metode au fost testate pentru probleme de clasificare netriviabile pentru care mulțimea inițială avea aproximativ 2000 exemple, implicând 14 atribute.

Rezultatele obținute au demonstrat că a fost posibilă construcția unui arbore de decizie, care clasifica toate exemplele din mulțime corect, el fiind obținut dintr-o fereastră care reprezenta doar o mică fracțiune a mulțimii inițiale. De asemenea, un lucru foarte important, este acela că timpul necesar pentru obținerea arborelui de clasificare crește liniar cu dificultatea problemei.

2. Utilizarea algoritmului ID3 pentru probleme cu caracter nedeterminist

Algoritmul ID3 prezentat a fost dezvoltat pentru rezolvarea problemelor deterministe de tip cauză-efect. Pornind de la o mulțime dată de exemple C , se construiește un arbore de decizie. Clasificarea unui nou exemplu implică un algoritm de parcurgere top-down. Se pornește din rădăcină și se urmează calea dictată de valorile de atribut existente în exemplul respectiv până se ajunge la o frunză, dacă exemplul poate fi clasificat cu acel arbore. În frunză există o singură clasă, adică un singur diagnostic cu grad de certitudine 1. Cu alte cuvinte, restricțiile de aplicare ale algoritmului de clasificare sunt următoarele:

- valorile atributelor unei căi trebuie să fie cunoscute;
- este posibil un singur diagnostic, pentru că o frunză conține o singură clasă.

Dar, de cele mai multe ori, problemele reale au un caracter nedeterminist. În cazul problemei de clasificare, nedeterminismul se reflectă în următoarele caracteristici:

- valorile atributelor nu sunt cunoscute în totalitate;
- pentru o cale din arbore corespund mai multe diagnostice-ce(clase) cu grade diferite de certitudine.

În acest caz, mulțimea inițială de exemple are următoarea formă:

p_1	p_2	· · · · ·	p_n	<i>diagnostic</i>
v_1^1	v_2^1	· · · · ·	v_n^k	d_1 cu grad de certitudine c_1
v_1^j	—	· · · · ·	v_n^l	d_2 cu grad de certitudine c_2

Figura 4. Tabel de exemple cu informație incompletă și incertă

unde $i, j, k < n$; liniuța orizontală reprezintă o valoare de atribut necunoscută.

În cazul unei valori de atribut necunoscute, se consideră că oricare dintre valorile permise ale atributului este posibilă și se generează exemplul respectiv pentru toate valorile atributului. Este posibil ca prin multiplicarea exemplului pentru diferite valori de atribut să se obțină exemple identice din punct de vedere al valorilor de atribut dar având clase diferite. Caz în care în arborele de clasificare va apare o frunză cu mai multe diagnostice cu grad de certitudine diferit. Figura 5 ilustrează acest caz.

p_1	p_2	p_3	p_4	<i>diagnostic</i>
v_1^1	v_2^1	v_3^2	v_4^3	d_1 cu grad de certitudine c_1
v_1^1	—	v_3^2	v_4^3	d_2 cu grad de certitudine c_2

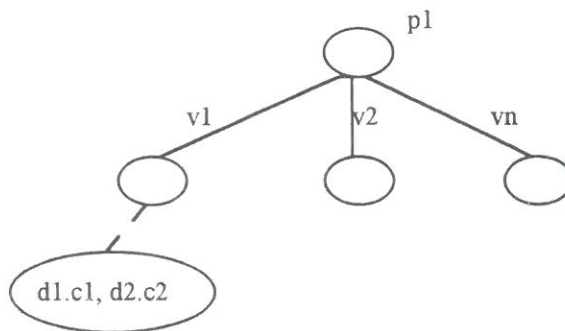


Figura 5. Exemple identice cu diagnostice diferite

2.1. Descrierea unei baze de cunoștințe ale cărei reguli au fost generate prin ID3

Așa cum a fost deja menționat, algoritmul ID3 poate fi utilizat pentru sintetizarea regulilor unei baze de cunoștințe, din exemple. Regulile sunt chiar arborii de clasificare obținuți dintr-o mulțime de exemple.

Sistemul cadru SCR (dezvoltat de autor, 1987) pentru generarea sistemelor expert bazate pe reguli, utilizează modelul cunoștințelor incomplete și incerte prezentat. Baza de cunoștințe este formată din reguli generate din exemple, utilizând algoritmul ID3, modificat pentru tratarea cazurilor de nedeterminism.

Structura unui modul al bazei de cunoștințe

Baza de cunoștințe a sistemului SCR este formată din mai multe module a căror structură este ilustrată în figura următoare:

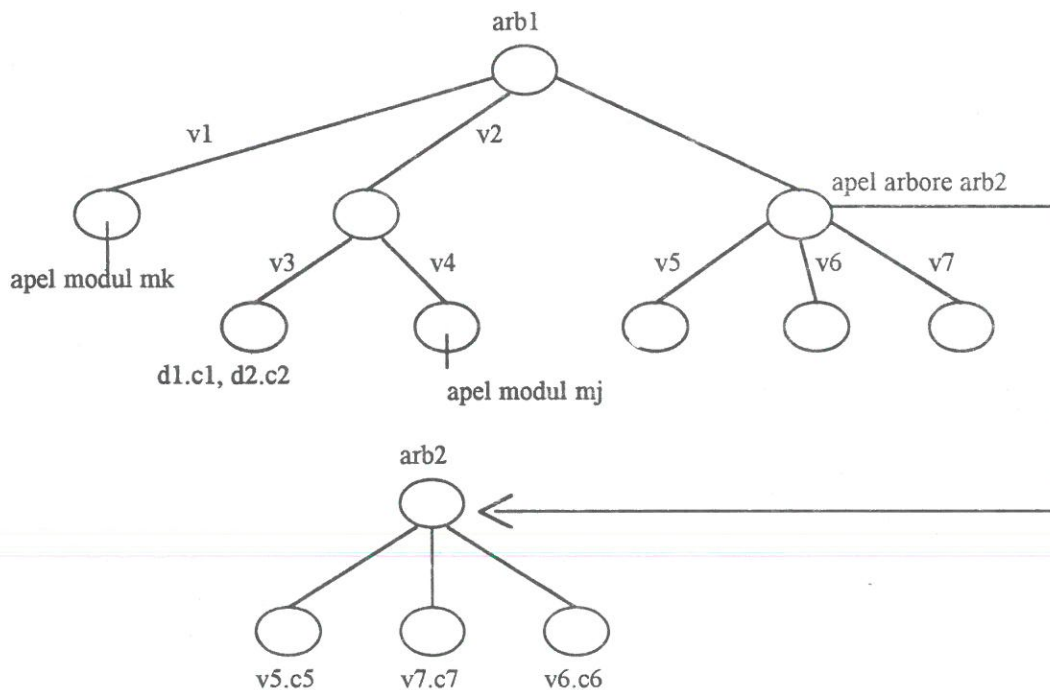


Figura 6. Structura unui modul din baza de cunoștințe

Un modul al bazei de cunoștințe poate fi format din mai mulți arbori care se caracterizează prin următoarele elemente particulare:

- un nod intern al unui arbore poate desemna o proprietate sau invocarea unui arbore, din același modul, pentru stabilirea valorilor atributului din nod; este cazul invocării arborelui *arb2*;
- o frunză poate conține unul sau mai multe diagnostice cu grade de certitudine diferite sau poate invoca un alt modul al bazei de cunoștințe, caz în care procesul de diagnoză continuă în acel modul.

Pentru a avea o reprezentare uniformă a unui arbore au fost impuse următoarele restricții de uniformitate:

- un modul poate conține un singur arbore;
- un nod poate invoca un alt modul:
- dacă nodul este o frunză atunci valorile întoarse de modul sunt diagnostice;

- dacă nodul este intern atunci valorile reîntoarse sunt folosite mai departe în procesul de diagnoză.

2.2. Controlul procesului de diagnoză

În cazul parametrilor cunoscuți a priori se selectează căile corespunzătoare valorilor respective. O dată ce calea este selectată, se realizează o "tăiere" a arborilor neselectați. Pentru atributele necunoscute care intervin, se evaluează atributul cel mai important, care trebuie obținut și se reia procesul.

Așa cum rezultă din structura unui modul al bazei de cunoștințe, valorile atributului unui nod care au fost obținute prin interogarea unui alt modul pot fi ponderate. Aceasta înseamnă că arcele arborelui pot avea ponderi și deci meritul unei frunze (deci al unei clase) depinde de ponderile căii care duce la frunză.

Această organizare a bazei de cunoștințe face posibilă ierarhizarea procesului de diagnoză, așa cum se observă și în figura următoare:

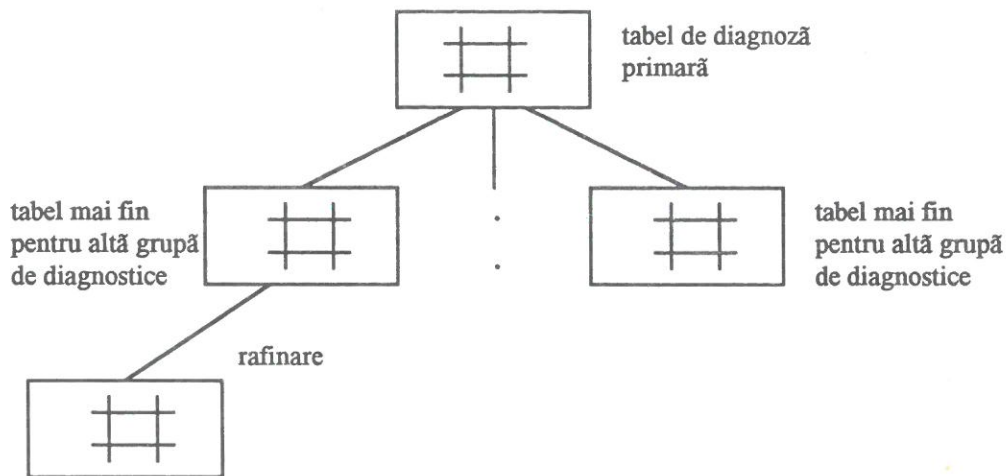


Figura 7. Rafinarea procesului de diagnoză

De exemplu, pentru diagnoza unui sistem de calcul, structura bazei de cunoștințe (modulele) ar putea fi următoarea:

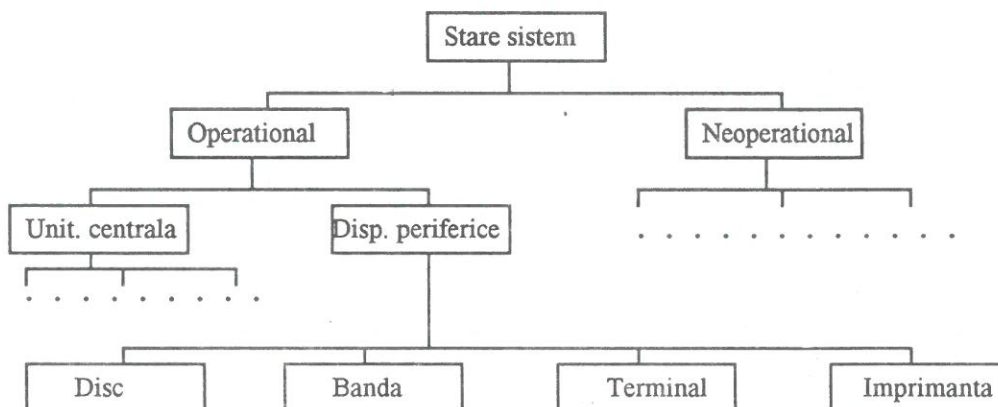


Figura 8. Diagnoza unui sistem de calcul

Avantajul cel mai important al ierarhizării procesului de diagnoză constă într-o rapidă diagnosticare prin restrângerea spațiului de căutare. Desigur există avantaje semnificative și din punct de vedere al construcției bazei de cunoștințe, un exemplu fiind paralelizarea construcției ei. De exemplu, introducerea exemplilor pentru modulul disc poate fi realizată în paralel cu oricare alt modul.

3. Optimizări ale algoritmului ID3

3.1. ID3 generalizat

Algoritmul ID3 realizează o căutare euristică, de tip "hill-climbing", în spațiul arborilor de decizie posibili. Căutările în spațiul stărilor de tip "hill-climbing" sunt căutări fără revenire care se bazează pe criterii de optim local în alegerea căii.

Dar nu întotdeauna optimele locale conduc la un optim global, astfel încât soluția determinată nu este întotdeauna cea mai bună. Din punct de vedere al algoritmului ID3, această "slăbiciune" constă în pierderea unor arbori de decizie mai buni, pentru aceleași date.

Există două aspecte ale unor arbori "mai puțin buni". Problema *valorilor irelevante* și problema *arcelor care lipsesc*. Atunci când ID3 alege un atribut pentru rădăcina subarborului curent, creează arce pentru toate valorile atributului respectiv care apar în exemple. Anumite valori pot fi relevante pentru clasificare, iar altele nu. Problema *valorilor irelevante* constă în crearea subarborilor pentru valori irelevante, care conduc la supraspecializarea regulilor. Se construiesc reguli pentru care trebuie verificate condiții care nu sunt necesare sau nu sunt relevante.

Problema *arcelor care lipsesc* constă, în esență, în reducerea capacității inductive a arborelui de clasificare. Ea este datorată faptului că

anumite submulțimi corespunzătoare unor noduri interne nu conțin exemple pentru toate valorile atributului din nod. O astfel de situație este ilustrată în figura 9.

A	B	Clasa
a_1	b_1	C_1
a_1	b_2	C_2
a_2	b_2	C_2
a_2	b_3	C_1
a_3	b_1	C_3
a_3	b_2	C_3

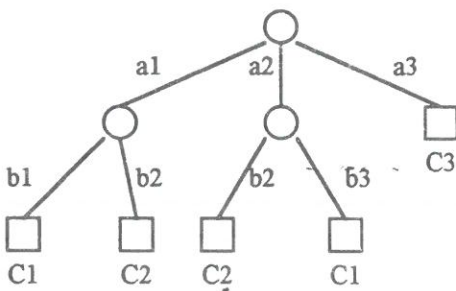


Figura 9. Arbore de decizie obținut prin ID3

Pentru exemplele :

$$e = A = a_3, B = b_3,$$

$$e' : A = a_1, B = b_3$$

arborele construit îl poate clasifica pe primul, dar nu și pe al doilea, pentru că în submulțimea corespunzătoare arcului etichetat cu a_1 nu există nici un exemplu pentru valoarea b_3 .

Pentru a rezolva acest tip de probleme Cheng&All a propus o versiune generalizată a algoritmului ID3 (GID3). În esență algoritmul este același, cu excepția faptului că nu se creează arce pentru toate valorile atributului. Se iau în considerare doar valorile semnificative, iar celelalte se grupează într-o valoare considerată implicită. Alegerea atributelor semnificative se face în funcție de un grad de toleranță. Valorile de atribut care ies din plaja de toleranță sunt grupate într-o valoare implicită. Prezentarea detaliată a algoritmului este în [1].

Arborele obținut aplicând algoritmul ID3 generalizat, este următorul:

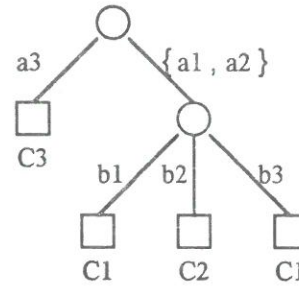


Figura 10. Arbore de decizie obținut prin GID3

Acest arbore poate clasifica ambele exemple.

Arborii generați utilizând algoritmul GID3 sunt mai compacți, mai fiabili și mai generali decât cei obținuți prin aplicarea algoritmului ID3. O calitate importantă a algoritmului GID3 este faptul că aceste îmbunătățiri sunt obținute fără o creștere semnificativă a complexității algoritmului.

3.2. ID3 incremental

Construcția arborilor de clasificare prin algoritmul ID3 este foarte bună în cazul în care baza de date (mulțimea de exemple) nu se modifică semnificativ. Pentru probleme în care noi exemple apar în mod regulat, ar fi de dorit ca arborele deja construit să poată fi modificat, în loc să fie reconstruit de fiecare dată.

Shlimmer și Fischer (1986) au creat algoritmul ID4 care construiește incremental arbori de clasificare, similari cu cei construiți de ID3. Ideea algoritmului este de elimina subarborii corespunzător unui atribut atunci când acesta trebuie înlocuit cu un altul mai bun.

Utgoff creează o variantă a algoritmului ID4, ID5, care diferă în modul de înlocuire a atributului de test. În loc să elimine subarborii corespunzător atributului care trebuie înlocuit, restructurează arborele propagând noul atribut. Propagarea atributului implică un proces de restructurare a arborelui fără pierdere de informație. El creează frunze pentru informația din exemple care nu se regăsește implicit în arborele existent astfel încât o propagare în sus să fie posibilă. Descrierea detaliată a algoritmului este prezentată în [4].

Costul utilizării algoritmului ID5 este mai mic decât cel al utilizării algoritmului ID3 pentru reconstrucția arborelui, la adăugarea unui nou exemplu care nu poate fi clasificat, cu excepția cazurilor în care atributele care trebuie înlocuite au un număr mare de valori.

Bibliografie

1. CHENG J. , USAMA, M. F., KEKI, B. I., ZHAOGANG, Q.: Improved Decision Trees: A Generalized Version of ID3. În: Proc. of the 5th Int. Conf. on MACHINE LEARNING, June 12-14, 1988.
2. HUNT, E. B., MARIN, J., STONE, P. T.: Experiments in Induction, Academic Press, New York, 1966.
3. QUINLAN, J. R.: Induction of Decision Trees. Machine Learning, Kluwer, 1986
4. UTGOFF, P. E.: ID5: An incremental ID3. În: Proc. of the 5th Int. Conf. on MACHINE LEARNING, June 12-14, 1988.