

Directii de Cercetare in Explorarea Datelor

Prof. Dan A. Simovici

University of Massachusetts Boston

Departmentul de Informatica

Boston, Massachusetts, USA

Ce este Explorarea Datelor (Data Mining)?

ED: Procesul de identificare a unor fapte si proprietati ale datelor.

- ED foloseste a varietate de discipline informatice: base de date, inteligenta artificiala, logica si statistica.
- ED este aplicata frecvent unor volume mari de date; exista numeroase grupuri de date de volum relativ redus care creaza probleme dificile.

Ce incercam sa descoperim cu ED?

Probleme majore:

- Descoperirea asocierilor intre obiecte.
- Gruparea obiectelor in multimi de obiecte similare. (*clustering*)
- Clasificarea obiectelor bazata pe proprietatile lor
- Evaluarea intersului faptelor si proprietatilor descoperite.
- Prepararea datelor (curatire, discretizare, etc.).

Cine are nevoie de ED?

- banci si cei care acorda credit;
- medici si biologi care incearca sa descopere cauzele bolilor si sa formuleze diagnostice;
- organizatii guvernamentale care incearca sa neutralizeze raufacatori;
- informaticieni care dirijeaza retele informatice si dezvoltata algoritmi pentru cercetarea internetului;
- ecologi si biologi intersati sa descopere surse de poluare,
- si multi altii...

Ce cunostinte practice sunt necesare?

- base de date relationale; SQL si folosirea lui in C++, Java, si alte limbaje;
- algoritmi care lucreaza cu o varietate de structuri de date;
- gestionarea depozitelor de date (data warehousing);
- cunosterea pachetelor de programe principale: Clementine, SAS, WEKA, etc.

Ce cunoștințe teoretice sunt necesare?

- diverse arii de matematica:
 - Clustering spații metrice
 - algebra liniară și analiză funcțională
 - Clasificare teoria informației
 - grafuri
 - Reguli de latici
 - asociere
- teoria complexității : NP- și #P-completitudine
- teoria informației;
- probabilități și statistică.

Baze de date si data mining

Tycho Brache (1546–1601)	colector de date multe date, dar n-a extras legile astronomice
Johannes Kepler (1571-1630)	minier de date

Clustering

Important for:

- condensarea datelor (prezentarea concisa a datelor);
- identificarea tendintelor in date.

A. K. Jain (1999): “nu exista un algorithm pentru clustering care este universal aplicabil in decoperirea oricarei structuri prezente in multimi multidimensionale de date”

Un exemplu de algorithm - clusterizarea incrementală

- date nominale
- clustering incremental

Caracteristica principală: Clustering-ul incremental formează grupuri adăugând în mod succesiv obiecte la grupuri (clusters) sau formând noi grupuri.

Date numerice si date nominale

- Date numerice:
 - inaltime: 1.82m, 1.25m, ...
 - temperatura: 38, 41, 54
- Date nominale:
 - culoare: rosu, verde, albastru,...
 - forma: patrat, romb, trapez, cerc,...

Distante pot fi definite in mod natural intre obiecte care au attribute numerice (folosind diferite metrice din R^n).

Dificultati cu datele nominale

Lipsa unei distante “naturale”: singura distanta ce se poate introduce este distanta Hamming, unde $d(o, o')$ este numarul de attribute in care o si o' sunt diferite.

Istoric

Algoritmi de grupare incrementală:

- Hartigan (1975)
- Fisher (1987) : COBWEB

Aplicatii

- F. Can et al.: baze de date de documente (1993–1995)
- Langford: detectarea focarelor de infectie din spitale (2001)
- J. Lin: serii temporale
- M. Charikar: regasirea dinamica a informatiei
- M. Ester: magazii de date (data warehouses)

Interesul clusteringului incremental

- Folosirea memoriei principale este minima
- Cerintele de timp cresc linear cu numarul de obiecte (scalable algorithm)

Sisteme de obiecte (SO)

Un sistem de obiecte este o pereche $\mathcal{S} = (S, H)$, unde

- S este o multime numita multimea de obiecte ale sistemului \mathcal{S} ,
- $H = \{A_1, \dots, A_m\}$ este o multime de functii definite pe S .

A_i (numit un atribut al lui \mathcal{S}) este o functie $A_i : S \longrightarrow E_i$, unde E_i este domeniul lui A_i .

Partitii

O partitie pe o multime S este o colectie nevida de parti ale lui S indexata de o multime I ,

$\pi = \{B_i | i \in I\}$ asa fel incit:

- $\bigcup_{i \in I} B_i = S$, si
- $i \neq j$ implica $B_i \cap B_j = \emptyset$.

B_i sunt *blocurile partitiei* π . Multimea partitiilor lui S este notata cu $\text{PART}(S)$.

Latticea Partitiilor

$\pi \leq \sigma$ daca fiecare block B al partitiei π este inclus intr-un block al partitiei σ .

Daca $\pi, \pi' \in \text{PART}(S)$ exista o partitie minimala π_1 astfel ca $\pi \leq \pi_1$ si $\pi' \leq \pi_1$; de asemenea, exista cea mai mare partitie π_2 pentru care $\pi_2 \leq \pi$ si $\pi_2 \leq \pi'$. Prima partitie se noteaza cu $\pi \vee \pi'$; a doua cu $\pi \wedge \pi'$.

Partitii generate de attribute

Un atribut A al sistemului $\mathcal{S} = (S, H)$ genereaza o partitie $\pi^A \in \text{parts}(S)$: doua obiecte apartin aceluiasi bloc al partitiei π^A daca au aceiasi proiectie pe A .

B_a^A : blocul lui π^A care consta din obiectele lui S care au componenta pe A egala cu a .

In baze de date relationale π^A se obtine folosind optiunea **group by** A al frazei **select** in standard SQL.

Partitii generate de multimi de attribute

	T		
tid	...	L	...
t_1	...	a_1	...
\vdots	\vdots	\vdots	\vdots
t_i	...	a_i	...
\vdots	\vdots	\vdots	\vdots
t_n	...	a_n	...

L generate o partitie
a multimii de obiecte

$t_j \equiv_L t_k$ daca si numai daca $t_j[L] = t_k[L]$

Notam cu π_L partitia
generata de L

Partitii si dependente func- tionale

$$T$$

tid	...	L	...	K	...
\vdots	\vdots	\vdots	\vdots	\vdots	\vdots
t_i	...	a_i	...	b_i	...
\vdots	\vdots	\vdots	\vdots	\vdots	\vdots
t_j	...	a_j	...	b_j	...
\vdots	\vdots	\vdots	\vdots	\vdots	\vdots
t_n	...	a_n	...	b_n	...

T satisface dependenta functionala $L \rightarrow K$ daca $a_i = a_j$ implica $b_i = b_j$ pentru i, j , adica, $t_i \equiv_L t_j$ implica $t_i \equiv_K t_j$, adica, $\pi_L \leq \pi_K$

Clusterizari ca partitii

O clusterizare a unui sistem de obiecte $\mathcal{S} = (S, H)$ este o partiție κ a multimii de obiecte S .

Scopul nostru: determinarea grupajelor κ pornind de la legaturile lor cu partițiile induse de atribute π^A .

Valuari si Metrici

- $v : \text{PART}(S) \longrightarrow \mathbb{R}$ definita de
 $v(\pi) = \sum_{i=1}^n |B_i|^2$, unde $\pi = \{B_1, \dots, B_n\}$ este o valoare inferioara pe $\text{PART}(S)$:

$$v(\pi \vee \sigma) + v(\pi \wedge \sigma) \geq v(\pi) + v(\sigma) \quad (1)$$

pentru $\pi, \sigma \in \text{PART}(S)$.

- Pentru fiecare valoare inferioara v , functia
 $d : (\text{PART}(S))^2 \longrightarrow \mathbb{R}$ data de
 $d(\pi, \sigma) = v(\pi) + v(\sigma) - 2v(\pi \wedge \sigma)$ este o distanta pe $\text{PART}(S)$.

Criteriul de Optimalitate

Se cauta o grupare $\kappa = \{C_1, \dots, C_n\} \in \text{PART}(S)$ astfel ca distanta totala de la κ la partiile atributelor:

$$D(\kappa) = \sum_{i=1}^n d(\kappa, \pi^{A_i})$$

sa fie **minima**.

Grupaje si partiile atributelor

$$d(\kappa, \pi^A) = \sum_{i=1}^n |C_i|^2 + \sum_{j=1}^{m_A} |B_{a_j}^A|^2 - 2 \sum_{i=1}^n \sum_{j=1}^{m_A} |C_i \cap B_{a_j}^A|^2,$$

AMICA

(**A** **M**etric **I**ncremental **C**lustering **A**lgorithm)

Fie $t \notin S$, si fie $Z = S \cup \{t\}$. Urmatoarele situatii pot avea loc:

1. obiectul t este adagat unui grup(cluster) existent C_k , sau
2. un nou grup, C_{n+1} este creat care consta doar din t .

Relativ la π^A , t se adauga blocului $B_{t[A]}^A$.

Obiectul se adauga unui cluster existent

$$\kappa_{(k)} = \{C_1, \dots, C_{k-1}, C_k \cup \{t\}, C_{k+1}, \dots, C_n\}$$

$$\pi^{A'} = \{B_{a_1}^A, \dots, B_{t[A]}^A \cup \{t\}, \dots, B_{a_{m_A}}^A\}$$

$$d(\kappa_{(k)}, \pi^{A'}) - d(\kappa, \pi^A) = 2|C_k \oplus B_{t[A]}^A|.$$

Cresterea minima a $d(\kappa_{(k)}, \pi^{A'})$ este data de:

$$\min_k \sum_A 2|C_k \oplus B_{t[A]}^A|.$$

Obiectul formeaza un nou cluster

$$\begin{aligned}\kappa' &= \{C_1, \dots, \dots, C_n, \{t\}\} \\ \pi^{A'} &= \{B_{a_1}^A, \dots, B_{t[A]}^A \cup \{t\}, \dots, B_{a_{m_A}}^A\}\end{aligned}$$

$$d(\kappa', \pi^{A'}) - d(\kappa, \pi^A) = 2|B_{t[A]}^A|.$$

Directie de actionare

$$D(\kappa') - D(\kappa) = \begin{cases} 2 \cdot \sum_A |C_k \oplus B_{t[A]}^A| & \text{in Case 1} \\ 2 \cdot \sum_A |B_{t[A]}^A| & \text{in Case 2.} \end{cases}$$

Daca $\min_k \sum_A |C_k \oplus B_{t[A]}^A| < \sum_A |B_{t[A]}^A|$ se adauga t la clusterul C_k pentru care $\sum_A |C_k \oplus B_{t[A]}^A|$ este minima; altfel se creaza un nou cluster cu un singur obiect.

Dificultatile grupajului incremental

- Algoritmii de grupare incrementală sunt afectate, în general, de ordinea de prelucrare a obiectelor.
- Fiecare algoritm procedează într-o manieră “hill-climbing” care produce minime locale (și nu globale).

Limitarea efectului ordonarii obiectelor

Am folosit tehnica “not-yet” introdusa de Roure si Talavera:

NOT-YET: Un nou grupaj este creat numai daca conditia

$$r(t) = \frac{\sum_A |B_{t[A]}^A|}{\min_k \sum_A |C_k \oplus B_{t[A]}^A|} < \alpha,$$

este satisfacuta, adica, numai daca effectul adaugarii obiectului t asupra distantei totale $r(t)$ este suficient de semnificativ.

$\alpha \leq 1$ este un parametru dat de utilizator (daca $\alpha = 1$ obiectele nu sunt trimise la buffer).

Algoritmul AMICA

Intrari: Setul de date S si α

Iesiri: clustering C_1, \dots, C_{nc}

Metoda:

```

nc = 0;  $\ell = 1$ ;
while  $S \neq \emptyset$  do
    select an object  $t$ ;  $S = S - \{t\}$ ;
    if  $\sum_A |B_{t[A]}^A| \leq \alpha \min_{1 \leq k \leq nc} \sum_A |C_k \oplus B_{t[A]}^A|$ 
        then
            nc ++; create a new single-object cluster  $C_{nc} = \{t\}$ ;
        else
             $r(t) = \sum_A |B_{t[A]}^A| / \min_{1 \leq k \leq nc} \sum_A |C_k \oplus B_{t[A]}^A|$ 
            if  $r(t) > 1$ 
                then  $k = \arg \min_k \sum_A |C_k \oplus B_{t[A]}^A|$ 
                    add  $t$  to cluster  $C_k$ ;
                else /* this means  $\alpha < r(t) \leq 1$  */
                    place  $t$  in NOT-YET buffer;
            end if;
        end if;
    end while;

```


Experimente cu date produse sintetic

- Date sintetice: produse de un algoritm care genereaza obiecte cu componente reale grupate in jurul unui numar dat de centre.
- Datele au fost discretizate following un numar specific de intervale de discretizare, ceea ce ne permite sa tratam datele ca date nominale.
- Am experimentat cu citeva multimi de date cu un numar crescind de obiecte, cu un numar crescind de dimensiuni, folosind citeva permutari ale obiectelor.
- Toate experimentele folosesc $\alpha = 0.95$.

Stabilitatea Grupurilor

- Experiment executat pe o baza de date care consta din 10,000 de obiecte (grupate in jurul a 6 centroizi)
- O prima aplicare a algoritmului genereaza 11 grupuri.
- Cele mai multe obiecte (9895) sunt concentrate in 6 grupuri, ceea ce reprezinta o buna aproximare a grupurilor “naturale” produse de algoritmul de generare.

AMICA este relativ imuna la permutari

Initial		Permutatare Aleatoare		
Cluster	Mar.	Cluster	Mar.	Distributie (cluster original)
1	1548	1	1692	1692 (2)
2	1693	2	1552	1548 (1), 3 (3), 1 (2)
3	1655	3	1672	1672 (5)
4	1711	4	1711	1711 (4)
5	1672	5	1652	1652 (3)
6	1616	6	1616	1616 (6)
7	1	7	85	85 (8)
8	85	8	10	10 (9)
9	10	9	8	8 (10)
10	8	10	1	1 (11)
11	1	11	1	1 (7)

Scalabilitate

Numar de obiecte	Timp pt. 3 permutari (ms)			Timp mediu (ms)
2000	131	140	154	141.7
5000	410	381	432	407.7
10000	782	761	831	794.7
20000	1103	1148	1061	1104

Setul de date CIUPERCI

- Setul de date contine 8124 descrieri de ciuperci si este tipic folosit pentru probleme de clasificare.
- Algoritmii de clasificare incearca sa determine daca un tip de ciuperca este comestibil sau otravitor.
- Atributul (otravitor/comestibil) este eliminat si AMICA a fost aplicat la setul de date fara acest atribut.

Rezultate experimentale

Cl. no.	O/C	Total	Procentul grupului dominant
1	825/2752	3577	76.9%
2	8/1050	1058	99.2%
3	1304/0	1304	100%
4	0/163	163	100%
5	1735/28	1763	98.4%
6	0/7	7	100%
7	0/192	192	100%
8	36/16	52	69%
9	8/0	8	100%

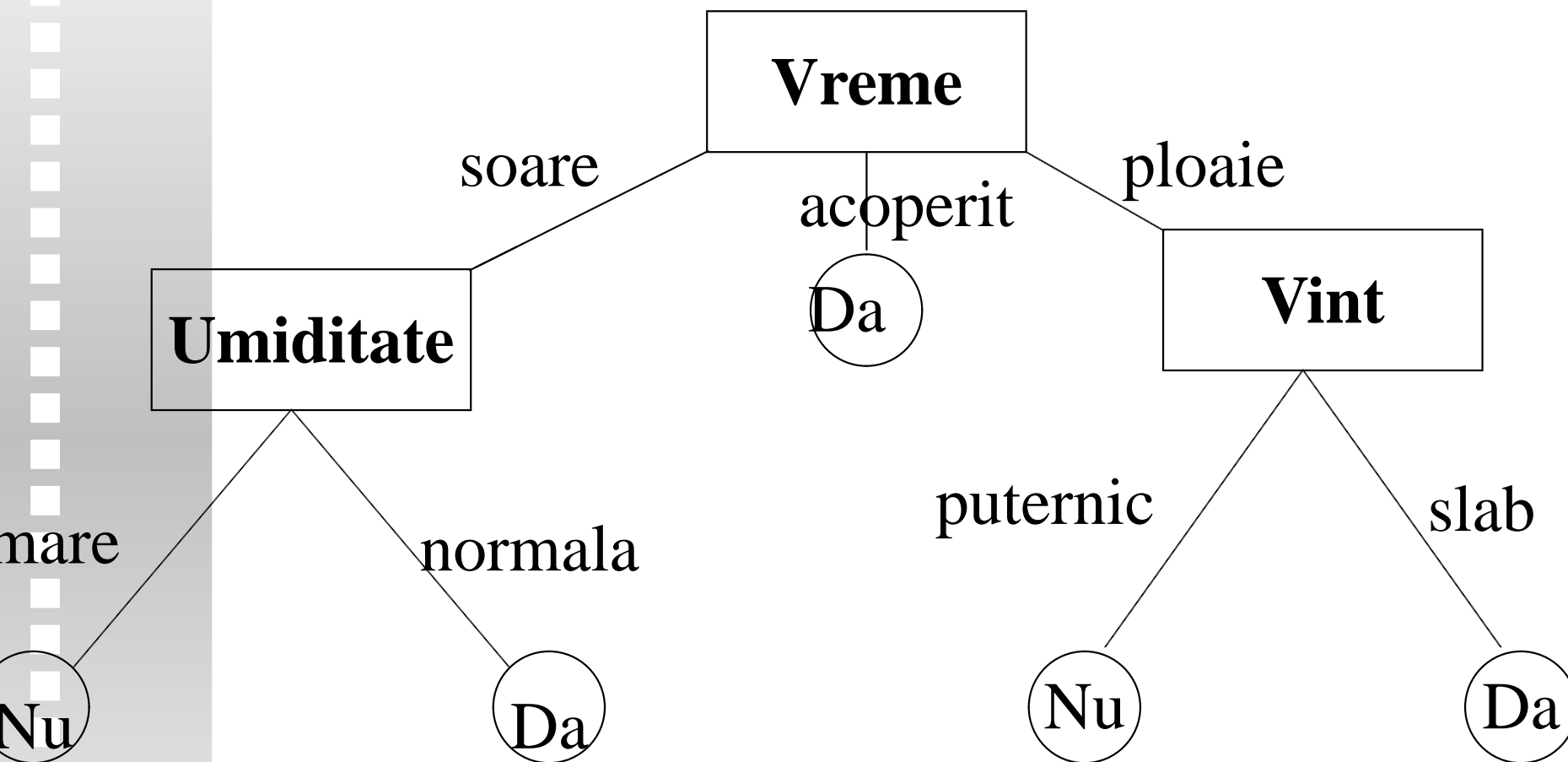
Stabilitate la Permutari

C_i	Grupuri Calculate									
	Permutare aleatoare									
	C'_1	C'_2	C'_3	C'_4	C'_5	C'_6	C'_7	C'_8	C'_9	C'_{10}
	3540	1797	1095	192	1296	8	36	7	137	16
3577	3540	0	37	0	0	0	0	0	0	0
1058	0	0	1058	0	0	0	0	0	0	0
1304	0	8	0	0	1296	0	0	0	0	0
163	0	26	0	0	0	0	0	0	137	0
1763	0	1763	0	0	0	0	0	0	0	0
7	0	0	0	0	0	0	0	7	0	0
192	0	0	0	192	0	0	0	0	0	0
52	0	0	0	0	0	0	36	0	0	16
8	0	0	0	0	0	8	0	0	0	0

Probleme inrudite

- Continuarea studiului experimental cu alte valori ale factorului “not-yet” α .
- Combinarea algoritmului AMICA cu tehnici speciale de discretizare pentru extinderea algoritmului la date cu caracter mix,
- Grupare incrementală în varianta “Semi-supervised” bazată pe AMICA.
- IC aplicat la date de tip “stream”

Arbori de decizie



Cum classifica arborii de decizie

(**Vreme** = soare, **Temperatura** = cald,
Umiditate = mare, **Vint** = puternic)

Orice arbore de decizie este reprezentat de o disjunctie de conjunctii:

$$\begin{aligned} & ((\mathbf{Vreme} = \text{soare} \wedge (\mathbf{Umiditate} = \text{normala})) \\ & \vee (\mathbf{Vreme} = \text{acoperit}) \\ & \vee ((\mathbf{Vreme} = \text{ploaie} \wedge (\mathbf{Vint} = \text{slab}))) \end{aligned}$$

Entropia lui Shannon

$$X : \begin{pmatrix} a_1 & \cdots & a_n \\ p_1 & \cdots & p_n \end{pmatrix}, \text{ where } p_1 + \cdots + p_n = 1.$$

Entropia lui X este

$$\mathcal{H}(X) = p_1 \log_2 \frac{1}{p_1} + \cdots + p_n \log_2 \frac{1}{p_n}.$$

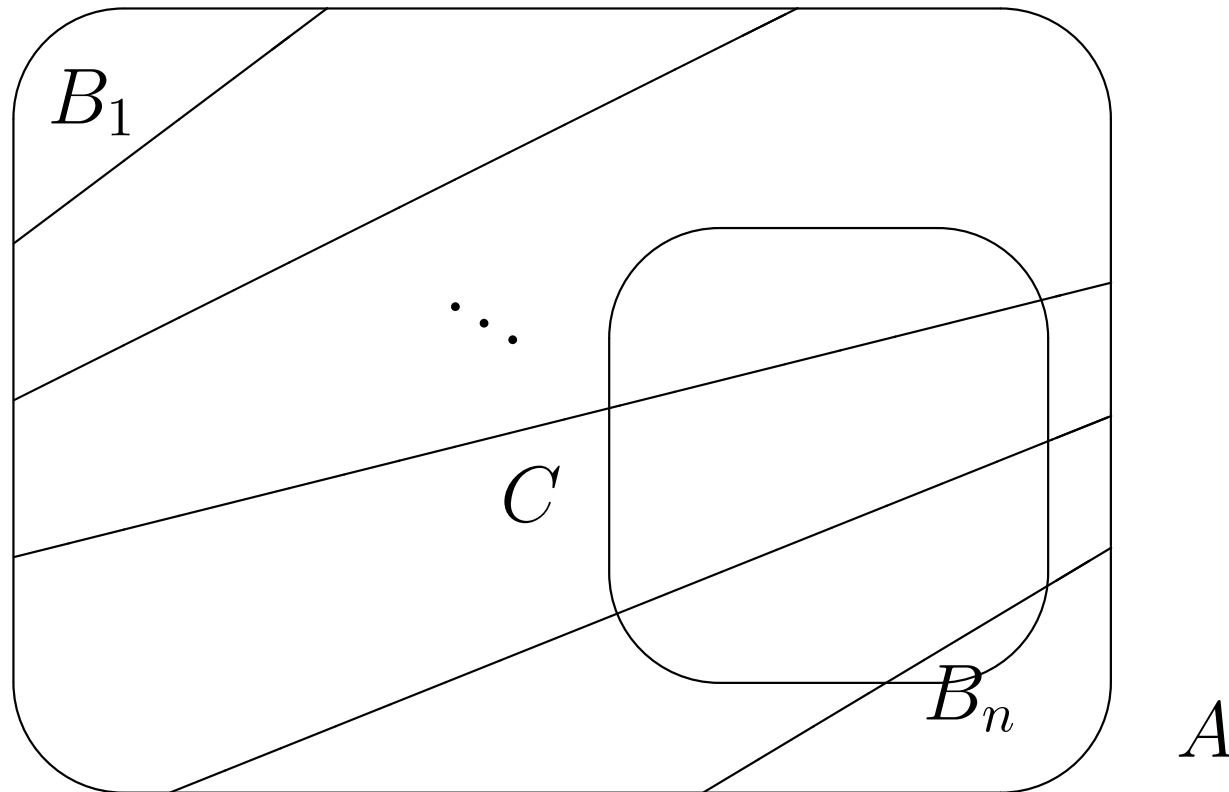
Daca $\pi = \{B_1, \dots, B_n\}$ este o partitie a multimii A atunci entropia lui π este:

$$\mathcal{H}(\pi) = - \sum_{i=1}^n \frac{|B_i|}{|A|} \log_2 \frac{|B_i|}{|A|}.$$

Urma unei partitii

Fie $\pi = \{B_1, \dots, B_n\}$ a partitie a multimii A si $C \subseteq A$.

Urma partitiei π pe C este $\pi_C = \{B_i \cap C \mid B_i \cap C \neq \emptyset\}$



Entropia Conditionala a Partitilor

$$\text{Fie } \begin{aligned} \pi &= \{B_1, \dots, B_n\} \\ \sigma &= \{C_1, \dots, C_m\} \end{aligned}$$

doua partitii ale multimii A . The **entropia conditionala** a lui π prin σ este:

$$\mathcal{H}(\pi|\sigma) = \sum_{j=1}^m \frac{|C_j|}{|C|} \mathcal{H}(\pi_{C_j})$$

Cistigul lui π relativ la σ este:

$$\text{Gain}(\pi, \sigma) = \mathcal{H}(\pi) - \mathcal{H}(\pi|\sigma)$$

Partitii si Arbori de Decisie

Alegera atributului de separare (splitting attribute) intr-un arbore de decizie se face in (ID3, sau C5.1 - Quinlan) folosind **cistigul informational**:

Fie K este atributul care defineste clasa, atunci alegerea atributului de separare A se face maximizind

$$\text{Gain}(\pi_K, \pi_A) = \mathcal{H}(\pi_K) - \mathcal{H}(\pi_K | \pi_A)$$

(Quinlan's ID3 or C4.5,...)

Zile favorabile pt. tenis

Zi	Vreme	Temp.	Umid.	Vint	Tenis
z1	soare	cald	rid	slab	nu
z2	soare	cald	rid	tare	nu
z3	acoperit	cald	rid	slab	da
z4	ploaie	mod	rid	slab	da
z5	ploaie	rece	nor	slab	da
z6	ploaie	rece	nor	tare	nu
z7	acoperit	rece	nor	tare	da
z8	soare	mod	rid	slab	nu
z9	soare	rece	nor	slab	da
z10	ploaie	mod	nor	slab	da
z11	soare	mod	nor	tare	da
z12	acoperit	mod	rid	tare	da
z13	acoperit	cald	nor	slab	da
z14	ploaie	mod	rid	tare	nu

$$\begin{aligned}\mathcal{H}(\pi_{tenis}) &= \\ & -\frac{5}{14} \log \frac{5}{14} \\ & -\frac{9}{14} \log \frac{9}{14} \\ & = 0.940\end{aligned}$$

Continuarea Exemplului

Penrtu vreme:

$$C_{soare} = \{z1, z2, z8, z9, z11\}$$

$$C_{acoperit} = \{z3, z7, z12, z13\}$$

$$C_{ploaie} = \{z4, z5, z6, z10, z14\}$$

Urmele partitiei π_{tenis} :

$$\pi_{tenis}C_{soare} = \{\{z1, z2, z8\}, \{z9, z11\}\}$$

$$\pi_{tenis}C_{acoperit} = \{\{z3, z7, z12, z13\}\}$$

$$\pi_{tenis}C_{ploaie} = \{\{z6, z14\}, \{z4, z5, z10\}\}$$

Urmele partitiei π_{tenis} :

$$\begin{aligned}\pi_{tenisC_{soare}} &= \{\{z1, z2, z8\}, \{z9, z11\}\} \\ \pi_{tenisC_{acoperit}} &= \{\{z3, z7, z12, z13\}\} \\ \pi_{tenisC_{ploaie}} &= \{\{z6, z14\}, \{z4, z5, z10\}\}\end{aligned}$$

$$\begin{aligned}\mathcal{H}(\pi_{tenisC_{soare}}) &= -\frac{3}{5} \log \frac{3}{5} - \frac{2}{5} \log \frac{2}{5} = 1.116 \\ \mathcal{H}(\pi_{tenisC_{acoperit}}) &= -\frac{4}{4} \log \frac{4}{4} = 0 \\ \mathcal{H}(\pi_{tenisC_{ploaie}}) &= -\frac{2}{5} \log \frac{2}{5} - \frac{3}{5} \log \frac{3}{5} = 1.116\end{aligned}$$

Calculule Similare

$$\text{Gain}(\pi_{tenis}, \pi_{vreme}) = 0.247$$

$$\text{Gain}(\pi_{tenis}, \pi_{umiditate}) = 0.151$$

$$\text{Gain}(\pi_{tenis}, \pi_{vint}) = 0.048$$

$$\text{Gain}(\pi_{tenis}, \pi_{vreme}) = 0.029$$

Atributul de scindare: **vreme**

Probleme generate de criteriul de cistig

- Alegerea atributului de scindare este pur locala. Arborele care rezulta nu este optimal in mod necesar.
- Arborii care rezulta pot avea multe virfuri terminale, ceea ce provoaca o fragmentare excesiva a datelor.

Metrici si arbori de decizie

López de Mántaras introduce o distanta bazata pe entropia Shannon.

$$d(\pi, \sigma) = \mathcal{H}(\pi|\sigma) + \mathcal{H}(\sigma|\pi).$$

Un nou criteriu de alegere a atributului de scindare:

$$A = \arg \min d(\pi_K, \pi_A)$$

Suma a doua partitii

Daca $M \cap P = \emptyset$ si

$$\pi = \{B_1, \dots, B_m\} \in \text{PART}(M),$$

$$\sigma = \{C_1, \dots, C_n\} \in \text{PART}(P),$$

definim $\pi + \sigma$ ca partitia multimii $M \cup P$:

$$\pi + \sigma = \{B_1, \dots, B_m, C_1, \dots, C_n\}.$$

Daca M, P, Q sunt disjuncte si

$\pi \in \text{PART}(M), \sigma \in \text{PART}(P), \tau \in \text{PART}(Q)$, atunci

$$\pi + (\sigma + \tau) = (\pi + \sigma) + \tau.$$

Axiomatizarea Entropiei Generalizate

Fie $\Phi : \mathbb{R}_{\geq 0}^2 \longrightarrow \mathbb{R}_{\geq 0}$ o functie continua, unde $\Phi(x, y) = \Phi(y, x)$, $\Phi(x, 0) = x$ pentru $x, y \in \mathbb{R}_{\geq 0}$ si $\beta \in \mathbb{R}$, $\beta > 0$.

Sistemul de axiome (Φ, β) pentru $\mathcal{H} : \text{PART}(A) \longrightarrow \mathbb{R}_{\geq 0}$ consta din

(P1) Daca $\pi, \pi' \in \text{PART}(A)$, $\pi \leq \pi'$, atunci $\mathcal{H}(\pi') \leq \mathcal{H}(\pi)$.

(P2) Daca A, B sunt doua multimi finite, $|A| \leq |B|$, atunci $\mathcal{H}(\iota_A) \leq \mathcal{H}(\iota_B)$.

(P3) Pentru $A, B, A \cap B = \emptyset, \pi \in \text{PART}(A)$ si $\sigma \in \text{PART}(B)$ avem:

$$\begin{aligned} \mathcal{H}(\pi + \sigma) \\ &= \left(\frac{|A|}{|A| + |B|} \right)^\beta \mathcal{H}(\pi) + \left(\frac{|B|}{|A| + |B|} \right)^\beta \mathcal{H}(\sigma) \\ &\quad + \mathcal{H}(\{A, B\}). \end{aligned}$$

(P4) Daca $\pi \in \text{PART}(A)$ si $\sigma \in \text{PART}(B)$, atunci

$$\mathcal{H}(\pi \times \sigma) = \Phi(\mathcal{H}(\pi), \mathcal{H}(\sigma)).$$

- β determina o entropie \mathcal{H}_β pina la un factor constant. β determina si functia Φ .
- Daca $\beta \neq 1$ atunci for a partition $\pi = \{A_1, \dots, A_n\} \in \text{PART}(A)$ we have:

$$\mathcal{H}_\beta(\pi) = \frac{k}{\beta - 1} \left(1 - \sum_{j=1}^n \left(\frac{|A_j|}{|A|} \right)^\beta \right),$$

unde k este o constanta astfel ca $k(\beta - 1) > 0$.

- Dacă $\beta \neq 1$ avem $\Phi(x, y) = x + y - \frac{1}{k}xy$ pentru $x, y \in \mathbb{R}_{\geq 0}$.
- Dacă $\beta = 2$ avem indexul Gini:

$$\mathcal{H}_2(\pi) = c \left(1 - \sum_{j=1}^n \left(\frac{|A_j|}{|A|} \right)^2 \right).$$

- Cazul limita $\beta \rightarrow 1$ da entropia Shannon, adica

$$\mathcal{H}_1(\pi) = -c \sum_{j=1}^n \frac{|A_j|}{|A|} \log_2 \frac{|A_j|}{|A|}.$$

si $\Phi(x, y) = x + y$ for $x, y \in \mathbb{R}_{\geq 0}$.

Entropia conditional data de (Φ, β) -entropy \mathcal{H} este
 $\mathcal{H}_\beta : \text{PART}^2 \longrightarrow \mathbb{R}_{\geq 0}$:

$$\mathcal{H}_\beta(\pi|\sigma) = \sum_{j=1}^n \left(\frac{|C_j|}{|A|} \right)^\beta \mathcal{H}_\beta(\pi_{C_j}),$$

unde $\pi, \sigma \in \text{PART}(A)$ si $\sigma = \{C_1, \dots, C_n\}$.
 $\mathcal{H}_\beta(\pi|\omega_A) = \mathcal{H}_\beta(\pi)$.

Daca $\pi \in \text{PART}(A)$ avem:

- $\mathcal{H}(\pi) = 0$ daca si numai daca $\pi = \omega_A$.
- Daca $\pi, \sigma \in \text{PART}(A)$ avem $\mathcal{H}_\beta(\pi|\sigma) = 0$ daca si numai daca $\sigma \leq \pi$.

- Fie $\pi, \sigma, \sigma' \in \text{PART}(A)$. Daca $\sigma \leq \sigma'$ atunci $\mathcal{H}_\beta(\pi|\sigma) \leq \mathcal{H}_\beta(\pi|\sigma')$ for $\beta > 0$.
- Fie $\pi, \sigma \in \text{PART}(A)$ si $\beta > 0$. Avem $\mathcal{H}_\beta(\pi|\sigma) \leq \mathcal{H}_\beta(\pi)$.
- Daca $\pi, \pi', \sigma \in \text{PART}(A)$ astfel ca $\pi \leq \pi'$ atunci $\mathcal{H}_\beta(\pi|\sigma) \geq \mathcal{H}_\beta(\pi'|\sigma)$.
- Pentru $\beta \geq 1$ avem $\mathcal{H}_\beta(\pi \wedge \sigma) \leq \mathcal{H}_\beta(\pi) + \mathcal{H}_\beta(\sigma)$.

Daca $\beta \geq 1$ si $\pi, \sigma, \tau \in \text{PART}(A)$ avem inegalitatea:

$$\mathcal{H}_\beta(\pi|\sigma) + \mathcal{H}_\beta(\sigma|\tau) \geq \mathcal{H}_\beta(\pi|\tau).$$

Rezultatul nostru generalizeaza rezultatul lui López de Mántaras:

Daca $\beta \geq 1$ fie $d_\beta : \text{PART}(A)^2 \longrightarrow \mathbb{R}_{\geq 0}$ definita de $d_\beta(\pi, \sigma) = \mathcal{H}_\beta(\pi|\sigma) + \mathcal{H}_\beta(\sigma|\pi)$ for $\pi, \sigma \in \text{PART}(A)$.
 d_β este o metrica pe $\text{PART}(A)$.

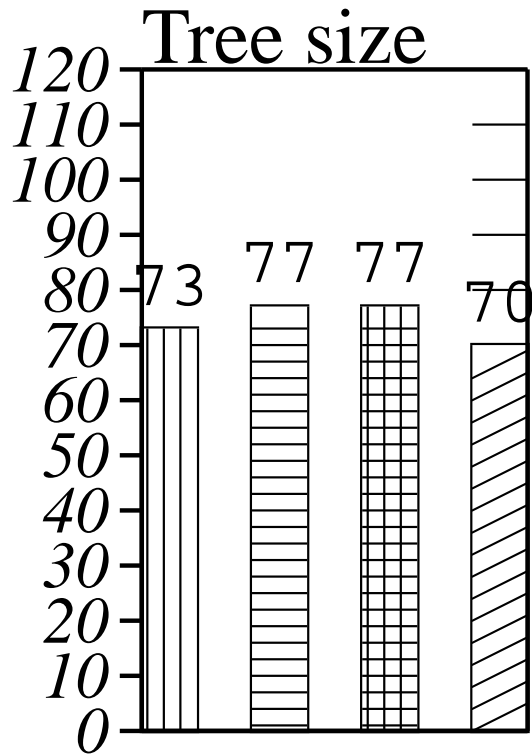
O noua alegere a atributului de scindare:

$$A = \arg \min d(\pi_K, \pi_A)$$

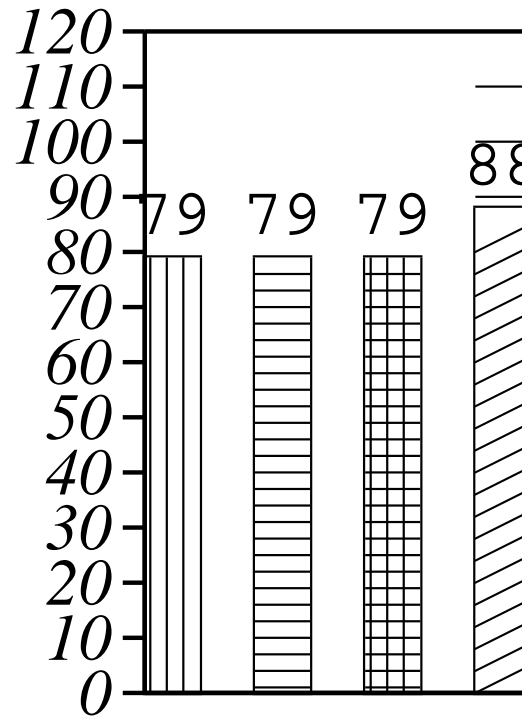
O noua problema: alegerea cea mai buna a parametrului β pentru o multime de date depinde de proprietatile ei statistice.

- Am experimentat cu 33 baze de date din colectia UCI.
- Fiecare experiment a folosit o 5-validare incrucisata; media a fost obtinuta pentru 5 experimente.

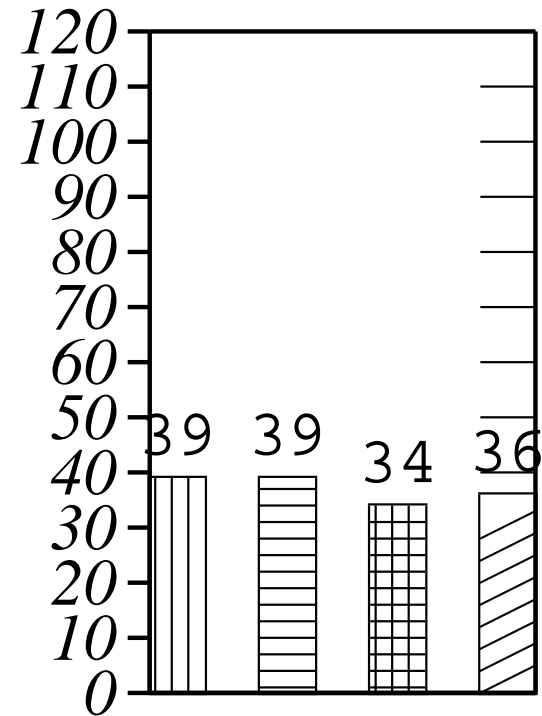
- Dimensiunea si numarul de virfuri terminale descreste pentru 18 din cele 33 baze de date si creste pentru celelalte 15.
- Cea mai importanta reducere a fost obtinuta pentru `primary-tumor`, unde numarul total de noduri a fost redus cu 37% pentru $\beta = 2.5$, iar numarul de noduri terminale a fost redus cu 38.8% in comparatie cu algoritmul standard (C5.0).



audiology

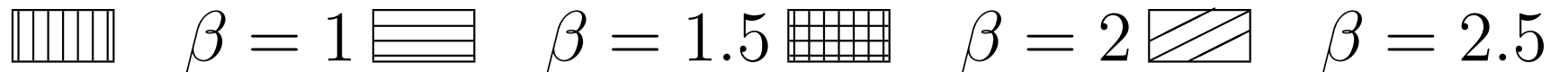


hepatitis

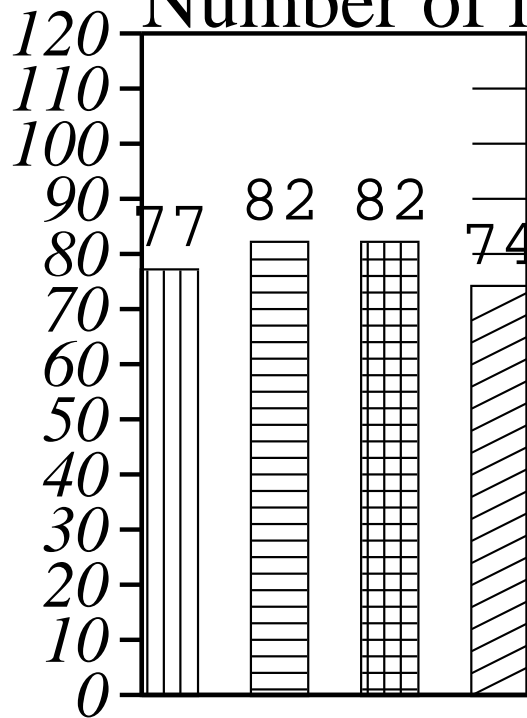


primary tumor

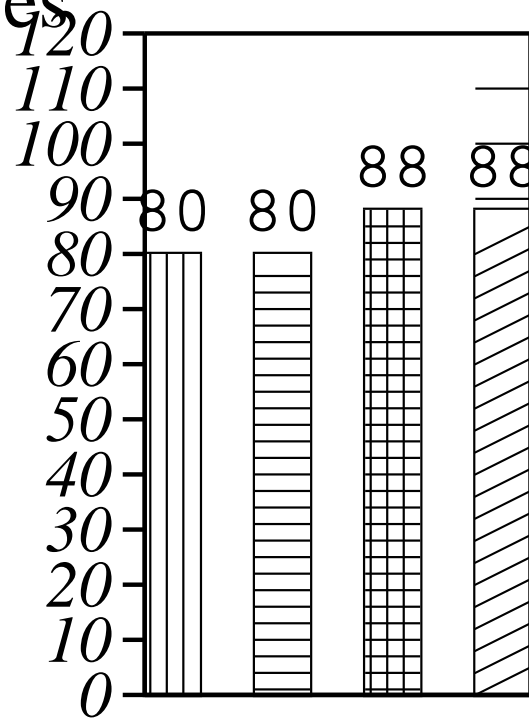
The β factor:



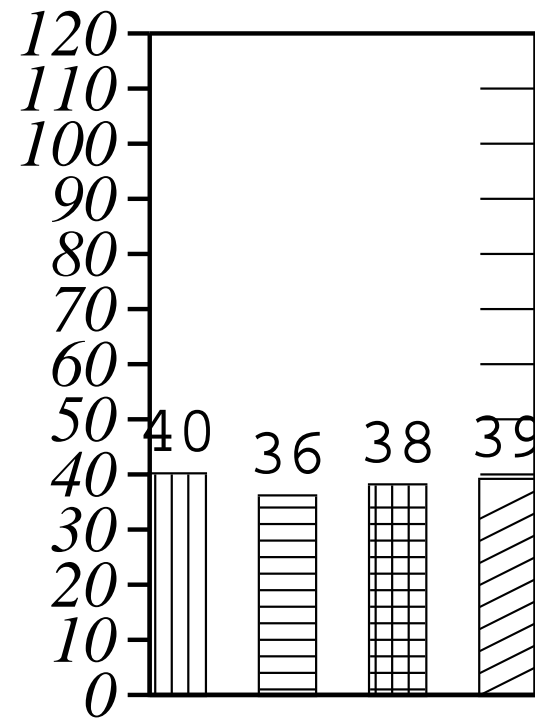
Number of leaves



audiology

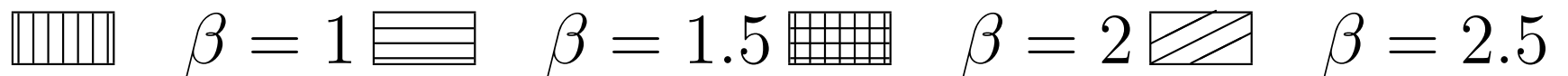


hepatitis



primary tumor

The β factor:



Unde ne putem informa despre DM?

- Conferinte principale:
 - KDD (USA)
 - PKDD (Europa)
 - PAKDD (Asia si Australia)
 - ICDM (anul acesta la Brighton, UK)
 - ICML
- TKDE (IEEE), Journal of Data Mining
- KDNuggets
- Internetul (CiteSeer)