

INTRODUCERE

Suntem copleșiți de date - date științifice, date medicale, date demografice, date financiare, date de marketing. Oamenii nu mai au timp să se uite la aceste date. Atenția umană a devenit o resursă importantă, astfel încât trebuiesc găsite căi de a analiza datele automat, de a le clasifica automat, de a le sintetiza automat, de a descoperi automat tendințe în date și de a caracteriza automat aceste tendințe. Acest "minerit în date" în vederea găsirii automate de cunoștințe și informații interesante/noi, este astăzi unul dintre cele mai active și interesante domenii de cercetare. Cercetătorii din domeniile bazelor de date, statisticii matematice, inteligenței artificiale și vizualizării computerizate sunt implicați și contribuie la dezvoltarea acestui domeniu.

Lucrarea de față prezintă tehnicile clasice "împrumutate" din statistica matematică de noul domeniu - am numit aici Data Mining; este vorba, mai precis, de tehnici de statistică exploratorie multidimensionale.

Statistica descriptivă permite reprezentarea vie și asimilabilă a informațiilor statistice prin simplificare și schematizare. Statistica descriptivă multidimensională este generalizarea naturală a cazului în care informațiile privesc mai multe variabile și/sau dimensiuni.

Trecerea la multidimensional implică însă o schimbare calitativă importantă. Într-adevăr, se spune despre microscop sau despre aparatul radiografic că nu sunt numai instrumente de descriere ci și instrumente de observație, de explorare și de cercetare. Prin metodele de statistică exploratorie multidimensională realitatea nu este doar simplificată pentru că este complexă, ci și explorată pentru că este ascunsă. Munca de pregătire și de codificare a datelor, regulile de interpretare și validare furnizate de tehnicile furnizate în cazul multidimensional, nu au simplitatea întâlnită în statistica descriptivă elementară. Nu este vorba doar de a prezenta ci și de a analiza, a descoperi, uneori de a verifica și dovedi, eventual de a testa anumite ipoteze.

Această lucrare conține prelegerile ținute studenților de la specializările INFORMATICĂ, și MATEMATICĂ APLICATĂ ale Facultății de Matematică și Informatică a Universității din București începând cu anul universitar 1995/1996, în cadrul unor cursuri opționale organizate anual sau semestrial în funcție de solicitări.

Numărul metodelor ce permit descrierea și explorarea tabelor rectangulare de date statistice (tabele de măsurători-observații, tabele de contingență, tabele de prezență-absență, sau tabele de incidență) este destul de mare. Metodele reținute pentru a fi prezentate au fost alese în funcție de posibilitățile pe care le au de a manipula tabele voluminoase, în funcție de transparența funcționării lor, în funcție de calitatea inserției în evantaiul metodelor ce sunt în mod real aplicabile și aplicate.

Două mari familii de metode răspund la aceste exigențe:

- [capitolul 1]: *metodele factoriale* bazate pe căutarea axelor principale (analiza în componente principale și analiza corespondențelor simple și multiple sunt metodele factoriale cele mai utilizate) care, produc în principal, vizualizări grafice plane sau spațiale ale obiectelor cercetate;
- [capitolul 2]: *metodele de clasificare* care produc agregări în clase de obiecte sau în familii de clase ierarhizate, obținute în urma unor calcule algoritmice. Obiectele cercetate sunt grupate, pornind de la vectorii care le descriu, în maiera cea mai puțin arbitrară.

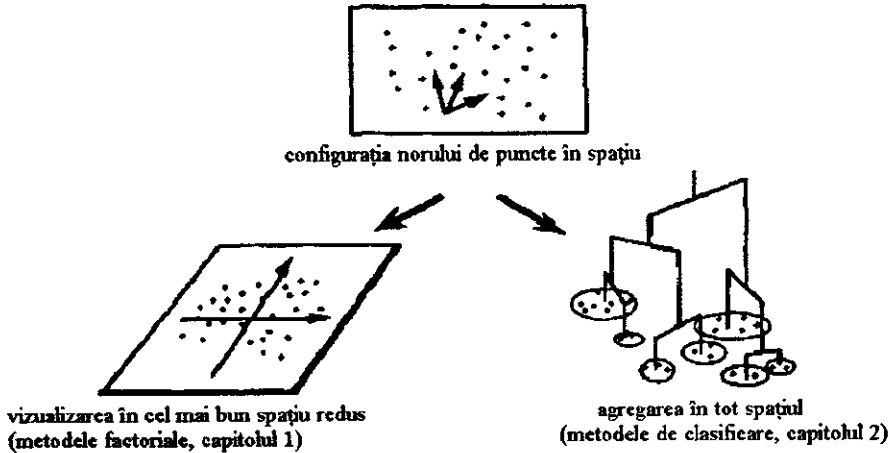


Figura 1 Cele două mari familii de metode ale statisticii exploratorii multidimensionale

Punctele de vedere furnizate de cele două tipuri de metode sunt în esență complementare. Vom insista asupra acestei complementarități care se manifestă de altfel la mai multe niveluri, fie că este vorba de posibilitatea de a înțelege structuri diverse, fie că este vorba de a ajuta lectura rezultatelor obținute.

- [capitolul 3]: *metodele explicative uzuale* vor lămuri pe utilizator asupra vocației specifice fiecărei metode (este vorba de analiza discriminantă și de metodele de segmentare) cât și asupra legăturilor cu metodele statisticii exploratorii (descrise în primele două capitole). Acest evantai de tehnici acoperă o parte importantă a aplicațiilor potențiale ale statisticii.

Nu există totuși o metodologie generală de articulare în practică a metodelor exploratorii de bază (metodele prezentate în capitolele 1 și 2) și metodele explicative uzuale (prezentate în capitolul 3). Fiecare aplicație implică, în funcție de domeniu și problemă, o muncă originală de codificare și selecție a metodelor particulare aplicate. În plus, trebuie să fim conștienți de faptul că metodele prezentate sunt eficiente în special în cazul datelor nestructurate sau amorfe (în care informația a priori asupra acestora este săracă).

Trebuie menționat faptul că există o literatură bogată privind tematica acestei lucrări. Bibliografia atașată constă numai dintr-o selecție a lucrărilor pe care autorul le-a consultat și care pot fi găsite cu ușurință în biblioteci.

Metodele prezentate au un pronunțat caracter matematic-aplicativ. Studenți, practicieni și cercetători din toate disciplinele ce trebuiesc să analizeze și să prelucereze volume mari de date multidimensionale, vor găsi în aceasta lucrare metodele de bază necesare.

Intenția autorului este de a continua dezvoltarea materialului prezentat aici într-o ediție următoare; în consecință observațiile și sugestiile sunt bine venite.

Autorul