

1 Ce este Data Mining?

Inițial “data mining” (extragerea de cunoștințe din date) a fost un termen din statistică însemnând suprautilizarea datelor pentru a deduce inferențe invalide.

- Teorema lui Bonferroni ne avertizează că atunci când există prea multe concluzii posibile unele dintre acestea vor fi adevărate din motive pur statistice, fără o valabilitate fizică.
- Exemplu faimos: David Rhine, un “parapsiholog” de la Duke a testat în anii 1950 studenții pentru “percepție extrasenzorială” (PE) cerându-le să ghicească 10 cărți de joc - roșii sau negre. A descoperit că 1/1000 din ei au ghicit toate cele 10 cărți și în loc să realizeze că este ceea ce trebuia așteptat din ghicirea aleatoare i-a declarat ca având PE. Când i-a retestat a descoperit că nu sunt mai buni decât media. Concluzia sa: a spune oamenilor că au PE duce la pierderea acesteia!

Definiția noastră: “descoperirea unor informații sumare utile despre date”.

1.1 Aplicații

Câteva exemple de “succes”:

1. Arbori de decizie construiți din istoria împrumuturilor bancare pentru a produce algoritmi de decizie în vederea acordării împrumuturilor.
2. Șabloane privind comportamentul călătorilor folosite pentru a gestiona vânzarea locurilor cu reducere la cursele aeriene, a camerelor de hotel, etc.
3. “Scutece și bere”: observația că cei care cumpără scutece cumpără bere mai mult decât media a permis supermarketurilor să plaseze berea și scutecele aproape unele de altele, știind că mulți cumpărători vor circula între ele. Plasarea cipsurilor de cartofi între cele două a crescut vânzările la cele trei articole.
4. Skycat și Sloan Sky Survey: gruparea obiectelor de pe cer după nivelul lor de radiație în diferite benzi a permis astronomilor să delimiteze galaxii, stele apropiate și alte feluri de obiecte cerești.
5. Compararea genotipului persoanelor îndeplinind sau nu o anumită condiție a permis descoperirea unei mulțimi de gene care împreună determină multe cazuri de diabet. Acest mod de extragere a cunoștințelor din date va deveni mult mai important în momentul construirii genomului uman.

1.1 Comunități implicate în data mining

Când data mining a fost recunoscută ca o unealtă puternică diferite comunități au clamat prioritate asupra subiectului:

1. Statistica
2. Inteligența artificială (IA) unde este denumită “machine learning”.
3. Cercetătorii din domeniul algoritmilor de grupare (“clustering algorithms”)
4. Cercetătorii din domeniul vizualizării datelor (“data visualization”)
5. Baze de date. Vom continua bineînțeles cu această abordare, concentrându-ne asupra provocărilor care apar atunci când cantitatea de date este mare iar calculele sunt complexe. Într-un anumit sens data mining poate fi văzută ca mulțimea algoritmilor pentru execuția unor cereri foarte complexe asupra unor date care nu sunt în memoria centrală a calculatorului.

1.2. Etape în procesul de extragere a cunoștințelor din date

1. *Colectarea datelor*, e.g. din depozitele de date (“data warehouse”), parcurgerea webului.
2. *Curățarea datelor*: eliminarea erorilor și/sau a datelor incorecte, e.g. temperatura pacientului=125
3. *Extragerea proprietăților*: obținerea doar a atributelor de interes ale datelor, e.g. “data achiziției” nu este probabil de interes în gruparea obiectelor cerești în cazul Skycat.
4. *Extragerea șablonelor și descoperirea*: Aceasta este etapa considerată adesea ca fiind “data mining” și aici ne vom concentra eforturile.
5. Vizualizarea datelor.
6. Evaluarea rezultatelor; nu orice fapt descoperit este și util ori adevărat!. O judecată este necesară înainte de a urma concluziile programelor de extragere de cunoștințe din date.