

## Capitolul 10

### Data mining – clustering

F. Radulescu, Curs: Utilizarea bazelor  
de date, anul IV CS.

1

### Problema

- ◆ Dându-se puncte într-un spațiu oarecare – deseori un spațiu cu foarte multe dimensiuni – grupează punctele într-un număr mic de *cluster*e, fiecare cluster constând din puncte care sunt "apropiate" într-un anumit sens.

F. Radulescu, Curs: Utilizarea bazelor  
de date, anul IV CS.

2

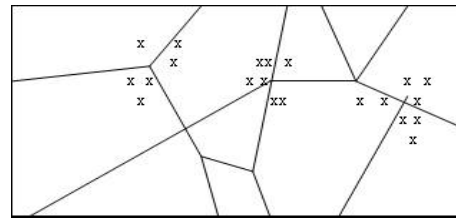
### Exemple de aplicatie

1. Cu mulți ani în urmă, în timpul unei izbucniri a holerei în Londra, un medic a marcat localizarea cazurilor pe o hartă, obținând un desen care arăta ca în figura următoare:

F. Radulescu, Curs: Utilizarea bazelor  
de date, anul IV CS.

3

### Exemple de aplicatie



F. Radulescu, Curs: Utilizarea bazelor  
de date, anul IV CS.

4

### Exemple de aplicatie

- ◆ Vizualizate corespunzător, datele au indicat că aparițiile cazurilor se grupează în jurul unor intersecții, unde existau puțuri infestate, arătând nu numai cauza holerei ci indicând și ce e de făcut pentru rezolvarea problemei.
- ◆ Din păcate nu toate problemele de data mining sunt atât de simple, deseori deoarece clusterelor sunt în atât de multe dimensiuni încât vizualizarea este foarte dificilă.

F. Radulescu, Curs: Utilizarea bazelor  
de date, anul IV CS.

5

### Exemple de aplicatie

2. *Skycat* a grupat în cluster 2 x 10<sup>9</sup> obiecte cerești în stele, galaxii, quasari, etc.
- ◆ Fiecare obiect era un punct într-un spațiu cu 7 dimensiuni, unde fiecare dimensiune reprezenta nivelul radiației într-o bandă a spectrului.
  - ◆ Proiectul Sloan Sky Survey este o încercare mult mai ambițioasă de a cataloga și grupa întregul univers vizibil.

F. Radulescu, Curs: Utilizarea bazelor  
de date, anul IV CS.

6

## Exemple de aplicatie

3. Documentele pot fi percepute ca puncte într-un spațiu multi-dimensional în care fiecare dimensiune corespunde unui cuvânt posibil.
- ◆ Poziția documentului într-o dimensiune este dată de numărul de ori în care cuvântul apare în document (sau doar 1 dacă apare, 0 dacă nu).
- ◆ Clusterelor de documente în acest spațiu corespund deseori cu grupuri de documente din același domeniu.

F. Radulescu, Curs: Utilizarea bazelor de date, anul IV CS.

7

## Distanța

- ◆ Pentru a discuta dacă o mulțime de puncte sunt suficient de apropiate pentru a fi considerate un cluster avem nevoie de o *măsură a distanței*  $D(x, y)$  care spune cât de depărtate sunt punctele  $x$  și  $y$ .
- ◆ Nu orice funcție poate fi utilizată ca funcție de măsurarea distanței.

F. Radulescu, Curs: Utilizarea bazelor de date, anul IV CS.

8

## Distanța

- ◆ Axiomele uzuale pentru o măsură a distanței  $D$  sunt următoarele:
1.  $D(x, y) \geq 0$
  2.  $D(x, x) = 0$ . Un punct este la distanță 0 de el însuși.
  3.  $D(x, y) = D(y, x)$ . Distanța e simetrică.
  4.  $D(x, y) \leq D(x, z) + D(z, y)$ . *Inegalitatea triunghiului.*

F. Radulescu, Curs: Utilizarea bazelor de date, anul IV CS.

9

## Distanța

- ◆ Deseori punctele pot fi percepute ca existând într-un spațiu euclidian  $k$ -dimensional și distanța între orice două puncte:
  - ◆  $x = [x_1, x_2, \dots, x_k]$  și
  - ◆  $y = [y_1, y_2, \dots, y_k]$
 este dată într-una din manierele uzuale:
  - ◆ Distanța comună ("norma L2")
  - ◆ Distanța *Manhattan* ("norma L1")
  - ◆ Maximul pe o dimensiune ("norma  $L_\infty$ ")

F. Radulescu, Curs: Utilizarea bazelor de date, anul IV CS.

10

## Distanța comună

- ◆ Distanța comună (sau norma L2) este dată de formula cunoscută:

$$\sqrt{\sum_{i=1}^k (x_i - y_i)^2}$$

F. Radulescu, Curs: Utilizarea bazelor de date, anul IV CS.

11

## Distanța Manhattan

- ◆ Distanța Manhattan (sau norma L1) este dată de formula următoare:

$$\sum_{i=1}^k |x_i - y_i|$$

- ◆ Ea poate fi folosită de exemplu și pentru calculul distanței între două puncte pe o placheta cu circuite imprimate multistrat.

F. Radulescu, Curs: Utilizarea bazelor de date, anul IV CS.

12

## Maximul pe o dimensiune

- ◆ Este data de formula:

$$\max_{i=1}^k |x_i - y_i|$$

- ◆ Aceasta funcție verifică toate cele 4 condiții pentru a fi funcție de distanță.
- ◆ Poate fi folosită de exemplu pentru spații euclidiene hiperdimensionale (număr de dimensiuni foarte mare)

F. Radulescu, Curs: Utilizarea bazelor de date, anul IV CS.

13

## Alte distante

- ◆ Unde nu există un spațiu euclidian în care să plasăm punctele gruparea devine mult mai dificilă.
- ◆ Iată un exemplu în care are sens: o măsură a distanței în lipsa unui spațiu euclidian

F. Radulescu, Curs: Utilizarea bazelor de date, anul IV CS.

14

## Alte distante

- ◆ Șirurile de caractere, cum sunt secvențele ADN, pot fi similare chiar și dacă există unele inserări și ștergeri precum și modificări ale unor caractere.
- ◆ De exemplu, *abcde* și *bcdxye* sunt destul de similare chiar dacă nu au nici o poziție comună și nu au nici chiar aceeași lungime.

F. Radulescu, Curs: Utilizarea bazelor de date, anul IV CS.

15

## Alte distante

- ◆ Astfel, în loc să construim un spațiu euclidian cu câte o dimensiune pentru fiecare poziție, putem defini funcția distanță:

$$D(x, y) = |x| + |y| - 2|LCS(x, y)|$$

unde LCS este cea mai lungă subsecvență comună lui  $x$  și  $y$ .

- ◆ În exemplul nostru  $LCS(abcde, bcdxye)$  este *bcd* de lungime 4, deci  $D(abcde, bcdxye) = 5 + 6 - 2 \times 4 = 3$ ; i.e. șirurile sunt destul de apropiate.

F. Radulescu, Curs: Utilizarea bazelor de date, anul IV CS.

16

## Alte distante

- ◆ Aceasta funcție de distanță arată câte caractere trebuie șterse sau adăugate unuia dintre șiruri pentru a obține celălalt șir.
- ◆ Într-adevăr, pentru a obține de exemplu pe *abcde* din *bcdxye* trebuie să:
  1. Adăugăm un *a* în față
  2. Ștergem pe *x*
  3. Ștergem pe *y*

F. Radulescu, Curs: Utilizarea bazelor de date, anul IV CS.

17

## Hiperdimensionalitatea

- ◆ O consecință mai puțin intuitivă a lucrului în spații hiperdimensionale este că aproape toate perechile de puncte sunt la o depărtare aproape egală cu media distanțelor între puncte.

Exemplu:

- ◆ Să presupunem că aruncăm aleator puncte într-un cub  $k$ -dimensional.
- ◆ Pentru  $k=2$ , ne așteptăm ca punctele să fie răspândite în plan cu unele foarte apropiate între ele și alte perechi aproape la distanță maxim posibilă.

F. Radulescu, Curs: Utilizarea bazelor de date, anul IV CS.

18

## Hiperdimensionalitatea

- ◆ Cu toate acestea, să presupunem  $k$  foarte mare, să zicem 100.000. Indiferent de norma folosită,  $L_2$ ,  $L_1$  sau  $L_\infty$ , știm că:

$$D(x, y) \geq \max_i |x_i - y_i|$$

pentru  $x = [x_1, x_2, \dots]$  și  $y = [y_1, y_2, \dots]$ .

- ◆ Pentru  $k$  foarte mare, e foarte posibil să existe o dimensiune  $i$  astfel încât  $x_i$  și  $y_i$  sunt diferite aproape de maximum posibil, chiar dacă  $x$  și  $y$  sunt foarte apropiate în alte dimensiuni.
- ◆ Astfel  $D(x, y)$  va fi foarte apropiată de 1.

F. Radulescu, Curs: Utilizarea bazelor de date, anul IV CS.

19

## Hiperdimensionalitatea

- ◆ O altă consecință interesantă a hiperdimensionalității este că toți vectorii,

$$x = [x_1, x_2, \dots] \text{ și } y = [y_1, y_2, \dots]$$

sunt aproape ortogonali.

- ◆ Motivul este că dacă proiectăm  $x$  și  $y$  pe oricare dintre cele  $k$  axe va exista unul în care proiecțiile vectorilor sunt aproape ortogonale (probabilitatea sa existe crește cu numărul de dimensiuni)

F. Radulescu, Curs: Utilizarea bazelor de date, anul IV CS.

20

## Abordari clustering

- ◆ La nivel înalt, putem împărți algoritmi de grupare în două mari clase:

1. Abordarea tip centroid: 'ghicim' centroizii sau punctele centrale pentru fiecare cluster și asignăm punctele la clusterul având cel mai apropiat centroid.

F. Radulescu, Curs: Utilizarea bazelor de date, anul IV CS.

21

## Abordari clustering

2. Abordarea ierarhică:

- ◆ Începem prin a considera că fiecare punct formează un cluster.
- ◆ Comparam repetat clusterurile apropiate prin folosirea unei măsuri pentru apropierea a două clusteruri (e.g. distanța dintre centroizii lor), sau pentru cât de bun va fi clusterul rezultat (e.g. distanța medie de la punctele din cluster la noul centroid rezultat).

F. Radulescu, Curs: Utilizarea bazelor de date, anul IV CS.

22

## Algoritmul k-means

- ◆ Acest algoritm este un algoritm popular care *ține datele în memoria centrală*
- ◆ Pe acest algoritm care se bazează alți algoritmi de clustering (ex.: BFR)
- ◆  $k$ -means alege  $k$  centroizi de cluster și asignează punctele la acestea alegând centroidul cel mai apropiat de punctul respectiv.
- ◆ Pe măsură ce punctele sunt asignare la un cluster, centroidul acestuia poate migra.

F. Radulescu, Curs: Utilizarea bazelor de date, anul IV CS.

23

## Observatie

- ◆ Centroidul unui cluster este 'centrul de greutate' al clusterului.
- ◆ El nu este de obicei unul din punctele care formează clusterul ci un punct între ele în spațiul euclidian respectiv.
- ◆ Conceptul de centroid nu este valabil decât în cazul clusteringului în spații euclidiene.
- ◆ Pentru cazurile în care nu avem un spațiu euclidian (punctele nu au coordonate) se folosește alternativ conceptul de clustroid.

F. Radulescu, Curs: Utilizarea bazelor de date, anul IV CS.

24

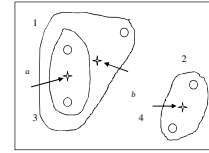
### Exemplu

- ◆ Un exemplu foarte simplu cu cinci puncte în două dimensiuni.
- ◆ Presupunem că asignăm punctele 1, 2, 3, 4 și 5 în această ordine, cu  $k=2$ .
- ◆ Atunci punctele 1 și 2 sunt asignate celor două clustere și devin centroidul lor pentru moment.

F. Radulescu, Curs: Utilizarea bazelor de date, anul IV CS.

25

### Exemplu



F. Radulescu, Curs: Utilizarea bazelor de date, anul IV CS.

26

### Exemplu

- ◆ Când considerăm punctul 3, să presupunem că este mai apropiat de 1, deci 3 se adaugă clusterului conținând 1 iar centroidul acestuia se mută în punctul marcat ca  $a$ .
- ◆ Presupunem că atunci când asignăm 4 găsim că 4 este mai aproape de 2 decât de  $a$ , deci 4 se alătură lui 2 în clusterul acestuia iar centrul se mută în  $b$ .
- ◆ În final, 5 este mai aproape de  $a$  decât de  $b$ , deci el se adaugă la clusterul  $\{1, 3\}$  al cărui centroid se mută în  $c$ .

F. Radulescu, Curs: Utilizarea bazelor de date, anul IV CS.

27

### Aplicare k-means

- ◆ Putem inițializa cei  $k$  centroizi alegând puncte suficient de depărtate de orice alt centroid până obținem  $k$ .
- ◆ Pe măsură ce calculul progresează putem decide să spargem un cluster și să unim două dintre ele pentru a păstra totalul de  $k$ . Testul pentru a decide asta poate fi să ne întrebăm dacă făcând operația respectivă se reduce distanța medie de la puncte la centroidul lor.

F. Radulescu, Curs: Utilizarea bazelor de date, anul IV CS.

28

### Aplicare k-means

- ◆ După localizarea centrozilor celor  $k$  clustere putem reasigna toate punctele deoarece unele puncte care au fost asignate la început pot acum să fie mai aproape de un alt centroid care s-a mutat.
- ◆ Dacă nu suntem siguri de valoarea lui  $k$  putem încerca valori diferite pentru  $k$  până când găsim cel mai mic  $k$  astfel încât mărirea lui  $k$  nu micșorează prea mult distanța medie a punctelor față de centroidul lor. Exemplul urmator ilustrează acest lucru.

F. Radulescu, Curs: Utilizarea bazelor de date, anul IV CS.

29

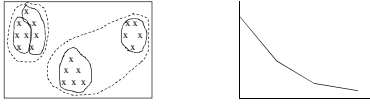
### Alt exemplu

- ◆ Să considerăm datele din figura următoare.
- ◆ În mod clar  $k=3$  este numărul corect de clustere dar să presupunem ca întâi încercăm  $k=1$ .
- ◆ În acest caz toate punctele sunt într-un singur cluster și distanța medie la centroid va fi mare.

F. Radulescu, Curs: Utilizarea bazelor de date, anul IV CS.

30

### Alt exemplu - cont



- ◆ Presupunem că apoi încercăm  $k=2$ .
- ◆ Unul dintre cele trei clusteruri va fi un cluster iar celelalte două vor fi forțate să creeze un singur cluster, așa cum arată linia punctată.
- ◆ Distanța medie a punctelor la centroid de va micșora astfel considerabil.

F. Radulescu, Curs: Utilizarea bazelor de date, anul IV CS.

31

### Alt exemplu - cont

- ◆ Dacă luăm  $k=3$  atunci fiecare dintre clusterurile vizibile va forma un cluster iar distanța medie de la puncte la centroizi se va micșora din nou, așa cum arată graficul din Figura.
- ◆ Totuși, dacă mărim  $k$  la 4 unul dintre adevăratele clusteruri va fi partiționat artificial în două clusteruri apropiate, așa cum arată liniile continue.

F. Radulescu, Curs: Utilizarea bazelor de date, anul IV CS.

32

### Alt exemplu - cont

- ◆ Distanța medie la centroid va scădea puțin dar nu mult.
- ◆ Acest eșec de a merge mai departe ne arată că valoarea  $k=3$  este corectă chiar dacă datele sunt în atât de multe dimensiuni încât nu putem vizualiza clusterurile.
- ◆ În acest fel aflăm valoarea corectă a lui  $k$  – numărul de clusteruri

F. Radulescu, Curs: Utilizarea bazelor de date, anul IV CS.

33

### Algoritmul BFR

- ◆ Bazat pe  $k$ -means, acest algoritm citește datele o singură dată în tranșe egale cu memoria centrală disponibilă la fiecare pas.
- ◆ Algoritmul lucrează cel mai bine dacă clusterurile sunt normal distribuite în jurul unui punct central, eventual cu o deviație standard diferită în fiecare dimensiune.

F. Radulescu, Curs: Utilizarea bazelor de date, anul IV CS.

34

### Reprezentare cluster în BFR

Cei care au creat acest algoritm și-au reprezentat clusterurile ca pe niște galaxii.

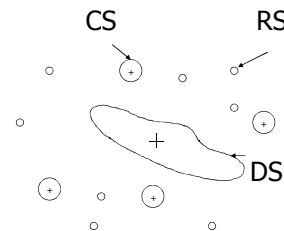
Un cluster constă din:

1. Un nucleu central numit **Discard set - DS**.
  2. Mulțimea acestor puncte este considerată ca aparținând în mod sigur clusterului.
- ◆ Toate punctele din această mulțime sunt înlocuite de niște statistici simple, descrise în continuare.
  - ◆ Notă: deși numite puncte « de aruncat » acestea au de fapt un efect semnificativ pe parcursul execuției algoritmului de vreme ce determină colectiv unde este centroidul și care este deviația standard a clusterului în fiecare dimensiune.

F. Radulescu, Curs: Utilizarea bazelor de date, anul IV CS.

35

### Reprezentare cluster în BFR



F. Radulescu, Curs: Utilizarea bazelor de date, anul IV CS.

36

## Reprezentare clustere în BFR

2. Galaxii înconjurătoare, numite colectiv **Compression set – CS** (*Mulțimea comprimată*).
- ◆ Fiecare subcluster din CS constă într-un grup de puncte care sunt suficient de apropiate unele de altele încât pot fi înlocuite cu statisticile lor, la fel ca și DS.
- ◆ Totuși, ele sunt suficient de departe de orice centroid de cluster încât nu suntem încă siguri de care cluster aparțin.

F. Radulescu, Curs: Utilizarea bazelor de date, anul IV CS.

37

## Reprezentare clustere în BFR

- ◆ Stele individuale care nu sunt parte a nici unei galaxii sau subgalaxii, **Mulțimea reținută (Retained set – RS)**.
- ◆ Aceste puncte nici nu pot fi asignate vreunui cluster nici grupate în vreun subcluster al CS.
- ◆ Ele sunt stocate în memoria centrală ca puncte individuale împreună cu statisticile pentru DS și CS.

F. Radulescu, Curs: Utilizarea bazelor de date, anul IV CS.

38

## Reprezentare comprimata

- ◆ Statisticile utilizate pentru a reprezenta fiecare cluster din DS și fiecare subcluster din CS sunt:

  1. Contorul numărului de puncte  $N$ .
  2. Vectorul sumelor coordonatelor punctelor în fiecare dimensiune. Vectorul este notat cu  $SUM_i$  iar componenta pentru dimensiunea  $i$  cu  $SUM_i$ .
  3. Vectorul sumelor pătratelor coordonatelor punctelor în fiecare dimensiune notat cu  $SUMSQ_i$ . Componenta pentru dimensiunea  $i$  cu  $SUMSQ_i$ .

F. Radulescu, Curs: Utilizarea bazelor de date, anul IV CS.

39

## Reprezentare - cont

- ◆ De notat că aceste trei tipuri de informații, totalizând în cazul în care avem  $k$  dimensiuni  $2k+1$  numere sunt suficiente pentru a calcula statistici importante pentru un cluster sau subcluster.
- ◆ Este mai convenabil de menținut pe măsură ce punctele sunt adăugate la cluster decât, să spunem, media și varianța în fiecare dimensiune.

F. Radulescu, Curs: Utilizarea bazelor de date, anul IV CS.

40

## Reprezentare - cont

- ◆ Cordonata  $\mu_i$  a centroidului clusterului în dimensiunea  $i$  este  $SUM_i/N$
- ◆ Varianța în dimensiunea  $i$  este:

$$\frac{SUMSQ_i}{N} - \left( \frac{SUM_i}{N} \right)^2$$

- ◆ iar deviația standard  $\sigma$  este rădăcina pătrată a acesteia.

F. Radulescu, Curs: Utilizarea bazelor de date, anul IV CS.

41

## Procesare

- ◆ La prima încărcare cu date a memoriei centrale, BFR selectează  $k$  centroizi de clustere utilizând un algoritm oarecare lucrând în memoria centrală, e.g. se ia un eșantion al datelor, se optimizează exact clusterelor și se aleg centroizii lor ca centroizi inițiali.
- ◆ O memorie centrală de puncte este procesată la fel în toate încărcările cu date urmatoare după cum urmează:

  1. Se determină care puncte sunt suficient de apropiate de un centroid curent astfel încât pot fi luate în DS iar statisticile lor ( $N$ ,  $SUM$ ,  $SUMSQ$ ) combinate cu statisticile anterioare ale clusterului.

F. Radulescu, Curs: Utilizarea bazelor de date, anul IV CS.

42

## Procesare

2. În memoria centrală se încearcă gruparea punctelor care nu au fost încă plasate în DS, inclusiv puncte ale RS din pașii precedenți.
- ◆ Dacă găsim un cluster de puncte a căror varianță este sub un prag ales, atunci vom privi aceste puncte ca un subcluster, le înlocuim cu statisticile lor și le considerăm parte a CS.
- ◆ Toate celelalte puncte vor fi plasate în RS.

F. Radulescu, Curs: Utilizarea bazelor de date, anul IV CS.

43

## Procesare

- ◆ Luăm în considerare unirea unui subcluster nou apărut cu un subcluster anterior din CS. Testul pentru a vedea dacă este de dorit să facem asta este ca mulțimea combinată de puncte să aibă o varianță sub un anumit prag. De notat că statisticile ținute pentru subclusterelor din CS sunt suficiente pentru a calcula varianța mulțimii combinate.

F. Radulescu, Curs: Utilizarea bazelor de date, anul IV CS.

44

## Procesare

- ◆ Dacă este ultimul pas, i.e. nu mai sunt date, atunci putem asigna subclusterelor din CS și punctele din RS la cel mai apropiat cluster de ele chiar dacă ele vor fi destul de departe de orice centroid de cluster.
- ◆ În felul acesta obținem clusterelor finale produse de algoritm

F. Radulescu, Curs: Utilizarea bazelor de date, anul IV CS.

45

## Scalarea multidimensională

- ◆ În multe cazuri nu avem un spațiu euclidian ci doar o mulțime de puncte și distanța între oricare două dintre acestea.
- ◆ Exemplu: un graf în care cunoaștem dimensiunea fiecărui arc. Din acestea putem afla distanța între oricare două noduri ca fiind lungimea drumului minim între ele

F. Radulescu, Curs: Utilizarea bazelor de date, anul IV CS.

46

## Scalarea multidimensională

- ◆ Se poate demonstra că având  $N$  puncte și distanțele între oricare 2 dintre ele putem crea un spațiu cu  $N-1$  dimensiuni
- ◆ În acest spațiu punctele sunt plasate exact (distanța calculată din coordonate este aceeași cu distanța de la care s-a plecat pentru orice pereche de puncte)

F. Radulescu, Curs: Utilizarea bazelor de date, anul IV CS.

47

## Scalarea multidimensională

- ◆ Problema este că în cazul unui număr mare de puncte rezultă un număr mare de dimensiuni (spațiu hiperdimensional)
- ◆ Ideal ar fi să plasăm cât mai exact cele  $N$  puncte într-un spațiu având  $K$  dimensiuni unde  $K \ll N$ .
- ◆ Acest proces se numește scalare multidimensională.
- ◆ Plasarea celor  $N$  puncte nu este 100% exactă (distanțele calculate din coordonatele rezultate nu sunt total exacte cu cele de la care s-a pornit)

F. Radulescu, Curs: Utilizarea bazelor de date, anul IV CS.

48



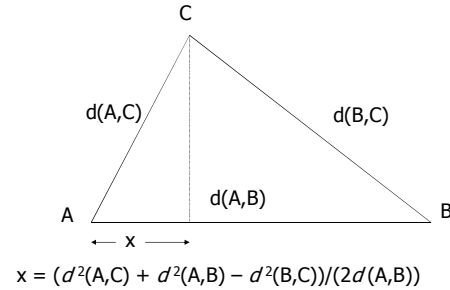
### Scalarea multidimensională

- ◆ Formula de baza folosită este cea prin care având două puncte putem afla distanțele proiecției unui al treilea punct pe segmentul format de primele două puncte.
- ◆ Formula este obținută din teorema lui Pitagora generalizată (teorema cosinusului)

F. Radulescu, Curs: Utilizarea bazelor de date, anul IV CS.

49

### Proiecția lui C pe AB



### Fastmap

- ◆ Algoritmul Fastmap este unul dintre algoritmii de scalare multidimensională.
- ◆ Acesta este un algoritm prin care se calculează succesiv coordonatele punctelor, câte o coordonată (o dimensiune) la fiecare pas.
- ◆ Pasul algoritmului este următorul:

F. Radulescu, Curs: Utilizarea bazelor de date, anul IV CS.

51

### Fastmap

1. Se aleg două puncte aflate la distanță cât mai mare, a și b. Acestea devin o axă de coordonate (cu originea în a).
2. Pentru orice punct c din cele N se calculează coordonata pe această axă conform formulei anterioare:

$$x = (D^2(a, c) + D^2(a, b) - D^2(b, c)) / (2 \cdot D(a, b))$$

F. Radulescu, Curs: Utilizarea bazelor de date, anul IV CS.

52

### Fastmap

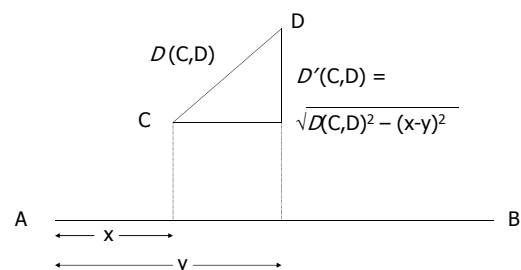
3. Pentru următoarele axe se vor folosi nu distanțele inițiale dintre puncte ci distanțe reportate în modul următor:  

$$D'^2 = D^2 - (x - y)^2$$
4. Procesul se sfârșește după calculul numărului dorit de coordonate sau când nu mai pot fi alese noi axe de coordonate.

F. Radulescu, Curs: Utilizarea bazelor de date, anul IV CS.

53

### Fastmap



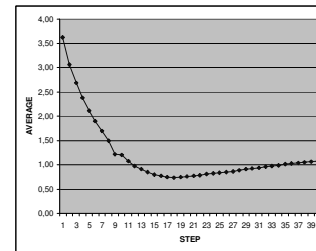
### Probleme in cazuri reale

- ◆ In cazurile reale (matricea nu este euclidiana) se poate intampla ca patrutul lui  $D'$  calculat cu formula anterioara sa dea un numar negativ.
- ◆ In astfel de cazuri pentru a putea continua se poate lua  $D' = 0$ .
- ◆ Aceasta alegere duce insa la erori care se propaga.
- ◆ Iata un exemplu de rulare pentru algoritmul in cazul in care sunt considerate 2000 de noduri.

F. Radulescu, Curs: Utilizarea bazelor de date, anul IV CS.

55

### Probleme in cazuri reale



F. Radulescu, Curs: Utilizarea bazelor de date, anul IV CS.

56

### Probleme in cazuri reale

- ◆ Figura arata media diferentei intre  $D_{real}$  si  $D_{calculat}$  unde:
- ◆  $D_{real}$  este distanta intre puncte de la care s-a pornit (cunoscuta prin ipoteza problemei)
- ◆  $D_{calculat}$  este distanta dintre puncte calculata pe baza coordonatelor obtinute pana la pasul respectiv.
- ◆ Se observa ca exista un minim dupa 18 pasi (spatiu optim are deci pentru acest exemplu 18 dimensiuni,  $k=18$ )

F. Radulescu, Curs: Utilizarea bazelor de date, anul IV CS.

57

### Probleme in cazuri reale

- ◆ In mod normal graficul ar trebui sa tinda asimptotic catre 0.
- ◆ Faptul ca nu se intampla asa e datorat erorii induse de considerarea lui  $D' = 0$  in cazul in care patrutul este negativ.
- ◆ Erorile acumulate fac ca dupa al 18-lea pas graficul sa inceapa sa creasca.

F. Radulescu, Curs: Utilizarea bazelor de date, anul IV CS.

58

### Bibliografie

- ◆ J.D.Ullman - CS345 --- Lecture Notes, Clustering I, II  
<http://infolab.stanford.edu/~ullman/cs345-notes.html>

F. Radulescu, Curs: Utilizarea bazelor de date, anul IV CS.

59

### Sfârșitul capitolului 10

F. Radulescu, Curs: Utilizarea bazelor de date, anul IV CS.

60