

Capitolul 9

Data mining – clustering

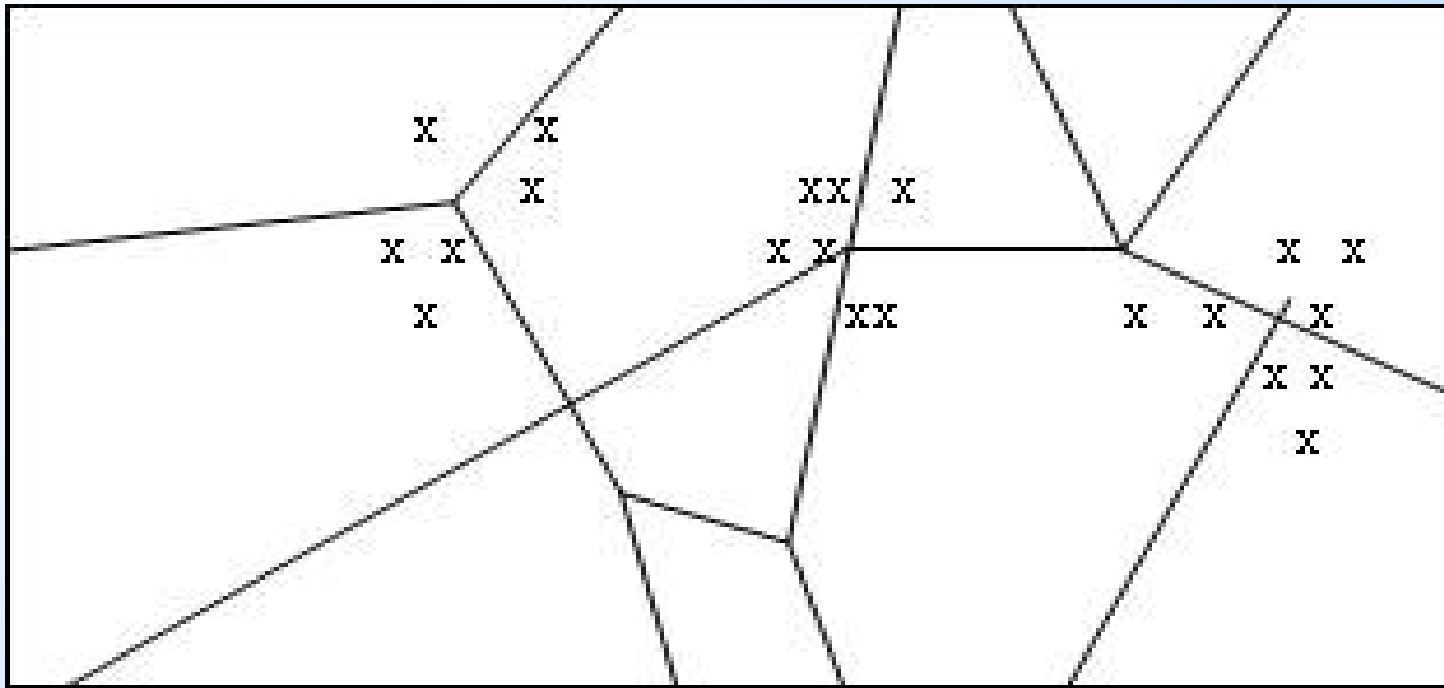
Problema

- ◆ 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.

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:

Exemple de aplicatie



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.

Exemple de aplicatie

2. *Skycat* a grupat în clusterne 2 x 10⁹ 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.

Exemple de aplicatie

3. Documentele pot fi percepute ca puncte într-un spațiu multi-dimensional în care fiecare dimensiune corespunde unui cuvânt posibil.
4. 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).
5. Clusterelor de documente în acest spațiu corespund deseori cu grupuri de documente din același domeniu.

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.

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.

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 L_2 ")
 - ◆ Distanța *Manhattan* ("norma L_1 ")
 - ◆ Maximul pe o dimensiune ("norma L_∞ ")

Distanța comună

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

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

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 imprimabile multistrat.

Maximul pe o dimensiune

- ◆ Este data de formula:

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

- ◆ Aceasta functie verifica toate cele 4 conditii pentru a fi functie de distanta.
- ◆ Poate fi folosita de exemplu pentru spatii euclidiene hiperdimensionale (numar de dimensiuni foarte mare)

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

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.

Alte distante

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

$$D(x, y) = |x| + |y| - 2|\text{LCS}(x, y)|$$

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

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

Alte distante

- ◆ Aceasta functie de distanta arata cate caractere trebuie sterse sau adaugate unuia dintre siruri pentru a obtine celalalt sir.
- ◆ Intr-adevar, pentru a obtine de exemplu pe *abcde* din *bcdxye* trebuie sa:
 1. Adaugam un a in fata
 2. Stergem pe x
 3. Stergem pe y

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 distanța maxim posibilă.

Hiperdimensionalitatea

- ◆ Cu toate acestea, să presupunem k foarte mare, să zicem 100.000. Indiferent de norma folosită, L_2 , L_1 sau L_∞ , 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 maximul posibil, chiar dacă x și y sunt foarte apropiate în alte dimensiuni.
- ◆ Astfel $D(x, y)$ va fi foarte apropiată de 1.

Hiperdimensionalitatea

- ◆ O alta consecință interesantă a hiperdimensionalității este că toți vectorii,
 $x = [x_1, x_2, \dots]$ și $y = [y_1, y_2, \dots]$
sunt aproape ortogonali.
- ◆ Motivul este că dacă proiectăm x și y pe oricare dintre cele k plane formate de două dintre cele k axe va exista unul în care proiecțiile vectorilor sunt aproape ortogonale (probabilitatea sa existe crește cu numărul de dimensiuni)

Abordari clustering

- ◆ La nivel înalt, putem împărți algoritmi de grupare în două mari clase:
 1. Abordarea tip centroid: 'ghicim' centriozii sau punctele centrale pentru fiecare cluster și asignăm punctele la clusterul având cel mai apropiat centroid.

Abordari clustering

2. Abordarea ierarhică:

- ◆ Începem prin a considera că fiecare punct formează un cluster.
- ◆ Comasăm repetat clusterelor apropiate prin folosirea unei măsuri pentru apropierea a două clusterelor (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).

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 asignate la un cluster, centroidul acestuia poate migra.

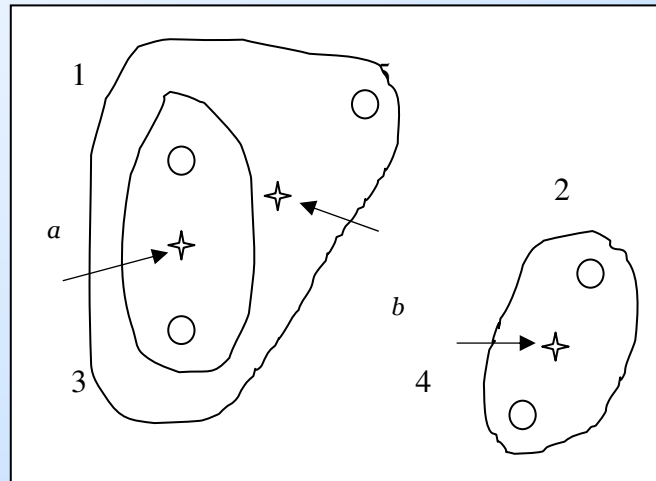
Observatie

- ◆ Centroidul unui cluster este 'centrul de greutate' al clusterului.
- ◆ El nu este de obicei unul din punctele care formeaza clusterul ci un punct intre ele in spatiul euclidian respectiv.
- ◆ Conceptul de centroid nu este valabil decat in cazul clusteringului in spatii euclidiene.
- ◆ Pentru cazurile in care nu avem un spatiu euclidian (punctele nu au coordonate) se foloseste alternativ conceptul de clustroid.

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 clusteruri și devin centroidul lor pentru moment.

Exemplu



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 .

Aplicare k-means

- ◆ Putem inițializa cei k centroizi alegând puncte suficient de depărtate de orice alt centroid până obținem k .
- ◆ Pe măsură ce calculul progresașază 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.

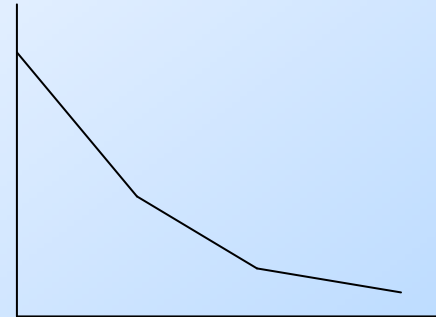
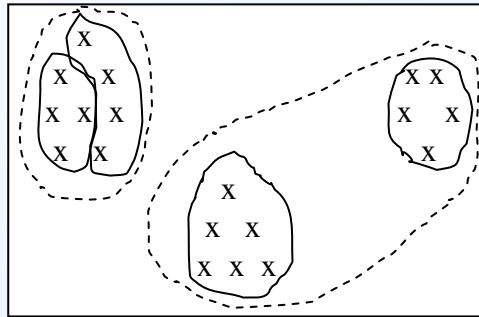
Aplicare k-means

- ◆ După localizarea centroizilor 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.

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.

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.

Alt exemplu - cont

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

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 clusterelor.
- ◆ In acest fel aflăm valoarea corectă a lui k – numărul de clusterelor

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 clusterelor sunt normal distribuite în jurul unui punct central, eventual cu o deviație standard diferită în fiecare dimensiune.

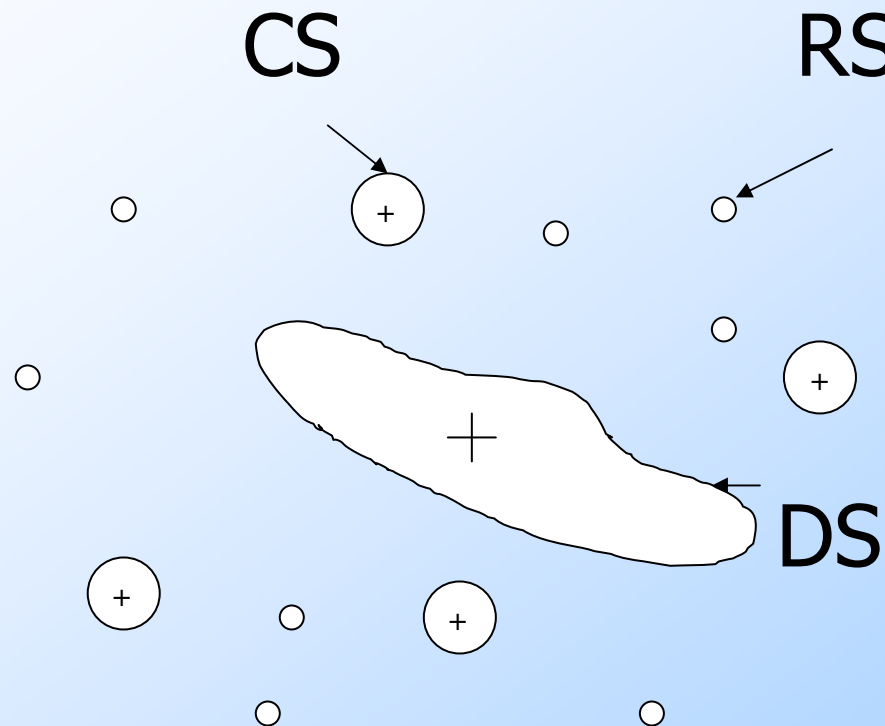
Reprezentare clusterare în BFR

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.

Reprezentare cluster în BFR



Reprezentare clusterare 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.

Reprezentare clusterare 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.

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 iar componenta pentru dimensiunea i cu SUM_i .
 3. Vectorul sumelor pătratelor coordonatelor punctelor în fiecare dimensiune notat cu $SUMSQ$. Componenta pentru dimensiunea i cu $SUMSQ_i$.

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.

Reprezentare - cont

- ◆ Cordonata μ_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 σ este rădăcina pătrată a acesteia.

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 următoare 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.

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.

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.

Procesare

- ◆ Dacă este ultimul pas, i.e. nu mai sunt date, atunci putem asigna subclustererele 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.
- ◆ In felul acesta obținem clustererele finale produse de algoritmul

Scalarea multidimensională

- ◆ In multe cazuri nu avem un spatiu euclidian ci doar o multime de puncte si distanta intre oricare doua dintre acestea.
- ◆ Exemplu: un graf in care cunoastem dimensiunea fiecarui arc. Din acestea putem afla distanta intre oricare doua noduri ca fiind lungimea drumului minim intre ele

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)

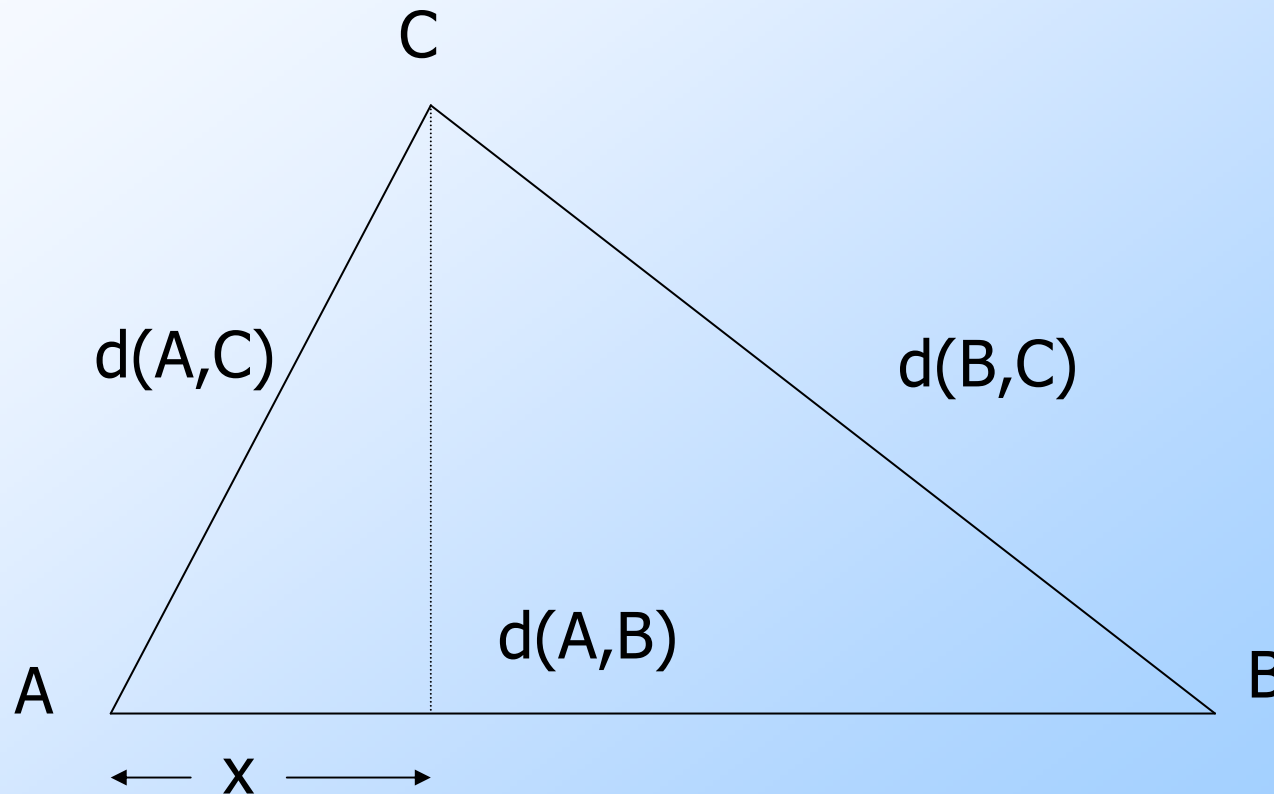
Scalarea multidimensională

- ◆ 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)

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)

Proiectia lui C pe AB



$$x = (d^2(A,C) + d^2(A,B) - d^2(B,C)) / (2d(A,B))$$

Fastmap

- ◆ Algoritmul Fastmap este unul dintre algoritmi 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:

Fastmap

1. Se aleg doua puncte aflate la distanta cat mai mare, a si b. Acestea devin o axa de coordonate (cu originea in a).
2. Pentru orice punct c din cele N se calculeaza coordonata pe aceasta axa conform formulei anterioare:

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

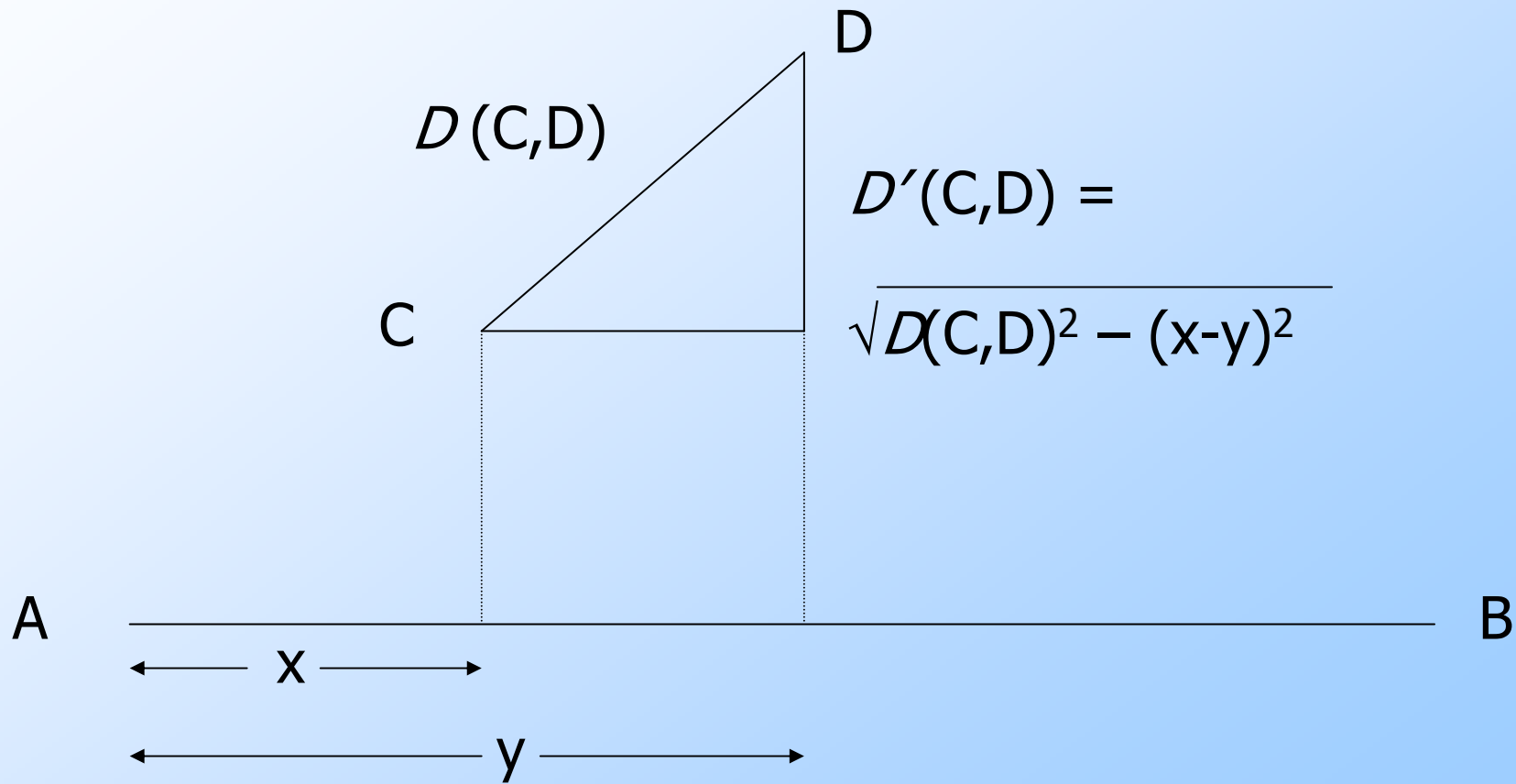
Fastmap

3. Pentru următoarele axe se vor folosi nu distanțele initiale 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.

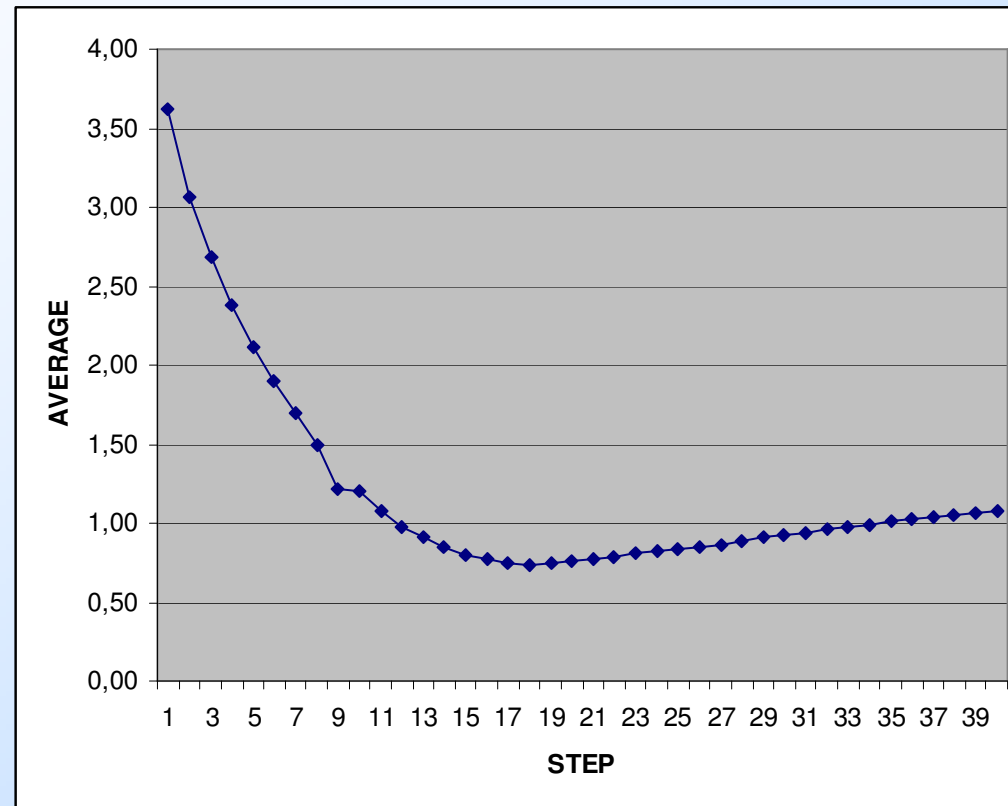
Fastmap



Probleme in cazuri reale

- ◆ In cazurile reale (matricea nu este euclidiană) 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.

Probleme in cazuri reale



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, $k=18$)

Probleme in cazuri reale

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

Bibliografie

- ◆ J.D.Ullman - CS345 --- Lecture Notes, Clustering I, II

<http://infolab.stanford.edu/~ullman/cs345-notes.html>

Sfârșitul capitolului 9