

## 7. Algoritmi de grupare (clustering)

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 "apropriate" într-un anumit sens. Câteva aplicații:

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 Fig. 14. 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 clusterurile sunt în atât de multe dimensiuni încât vizualizarea este foarte dificilă.

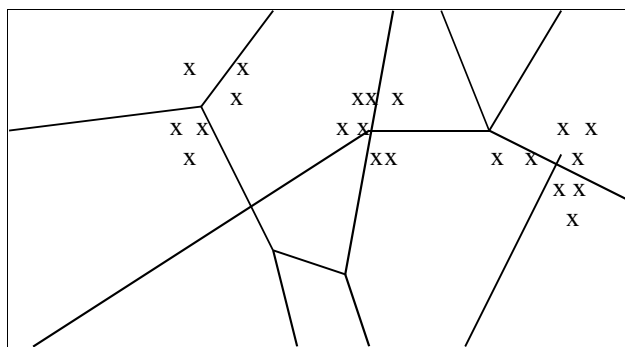


Figura 14: Clustere ale cazurilor de holeră au indicat unde erau puțurile infestate.

2. *Skycat* a grupat în clusteruri  $2 \times 10^9$  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.
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). Clusterurile de documente în acest spațiu corespund deseori cu grupuri de documente din același domeniu.

### 7.1 Măsuri ale distanței

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  $y$ . Axiomele uzuale pentru o măsură a distanței  $D$  sunt:

1.  $D(x, x) = 0$ . Un punct este la distanță 0 de el însuși.
2.  $D(x, y) = D(y, x)$ . Distanța e simetrică.
3.  $D(x, y) \leq D(x, z) + D(z, y)$ . **Inegalitatea triunghiului.**

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:

1. Distanța comună ("norma  $L_2$ "):  $\sqrt{\sum_{i=1}^k (x_i - y_i)^2}$ .
2. Distanța *Manhattan* ("norma  $L_1$ "):  $\sum_{i=1}^k |x_i - y_i|$ .
3. Maximul pe o dimensiune ("norma  $L_\infty$ "):  $\max_{i=1}^k |x_i - y_i|$ .

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

**Exemplul 7.1:** Putem privi paginile web ca puncte într-un spațiu generic  $10^8$ -dimensional unde fiecare dimensiune corespunde unui cuvânt. Totuși, calcularea distanțelor ar fi prohibitivă într-un spațiu atât de mare. O abordare mai bună este să legăm distanța  $D(x, y)$  de produsul scalar al vectorilor corespunzând lui  $x$  și  $y$ , din moment ce apoi vom lucra doar cu cuvinte care se află atât în  $x$  cât și în  $y$ .

Trebuie să calculăm lungimile vectorilor implicați și acesta sunt egale cu rădăcinile pătrate ale sumelor pătrătelor numerelor de apariție ale fiecărui cuvânt. Suma produselor numerelor de apariție ale fiecărui cuvânt în fiecare document se împarte la produsul lungimilor pentru a obține un produs scalar normalizat. Scădem această cantitate din 1 pentru a afla distanța dintre  $x$  și  $y$ .

De exemplu, să presupunem că sunt doar patru cuvinte de interes și vectorii  $x = [2, 0, 3, 1]$  și  $y = [5, 3, 2, 0]$  reprezintă numărul de apariții al acestora în două documente. Produsul scalar este  $2 \times 5 + 0 \times 3 + 3 \times 2 + 1 \times 0 = 16$ , lungimea primului vector este  $\sqrt{2^2 + 0^2 + 3^2 + 1^2} = \sqrt{14}$ , iar lungimea celui de-al doilea este  $\sqrt{5^2 + 3^2 + 2^2 + 0^2} = \sqrt{38}$ . Astfel,

$$D(x, y) = 1 - \frac{16}{\sqrt{14}\sqrt{38}} = 0.304$$

Ca un alt exemplu, să presupunem că  $x$  are vectorul de cuvinte  $[a_1, a_2, \dots]$  și  $y$  conține două copii ale lui  $x$ ; i.e.,  $y = [2a_1, 2a_2, \dots]$ . Atunci,

$$D(x, y) = 1 - \frac{\sum_i 2a_i^2}{\sqrt{\sum_i a_i^2} \sqrt{\sum_i (2a_i)^2}} = 0$$

Asta înseamnă că  $x$  și  $y$  sunt esențialmente același document; în mod sigur sunt despre același subiect și vor fi grupate în același cluster.

**Exemplul 7.2:** 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. 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.

## 7.2 Blestemul dimensionalității

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.

**Exemplul 7.3:** 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,  $\sqrt{2}$ .

Cu toate acestea, să presupunem  $k$  foarte mare, să zicem 100 sau 100.000. Indiferent de norma folosită,  $L_2$ ,  $L_1$  sau  $L_\infty$ , știm că  $D(x, y) \geq \max_i |x_i - y_i|$ , dacă  $x = [x_1, x_2, \dots]$  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.

O alta consecință interesantă a hiper-dimensionalității 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  $C_k^2$  plane formate de două dintre cele  $k$  axe va exista unul în care proiecțiile vectorilor sunt aproape ortogonale

## 7.3 Abordări pentru grupare (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.

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).

## 7.4 Algoritmi prezentați

Vom considera următorii algoritmi de grupare; ei se diferențiază după cum utilizează sau nu o distanță euclidiană și după abordarea folosită, de tip centroid sau ierarhică.

1. BFR [Bradley-Fayyad-Reina]: De tip centroid; utilizează o măsură euclidiană, cu clusterelor formate în jurul centroidului printr-un proces gaussian în fiecare dimensiune.
2. Fastmap: Nu e realment un algoritm de grupare ci un mod de a construi un spațiu euclidian cu puține dimensiuni pornind de la o măsură a distanței.
3. GRGPF: [Ganti et al.]: De tip centroid dar utilizează doar o măsură a distanței și nu un spațiu euclidian.
4. CURE: Ierarhic și euclidian, acest algoritm se ocupă de clusterelor având forme neobișnuite.

## 7.5 Algoritmul k-means

Acest algoritm este un algoritm popular care ține datele în memoria centrală și pe care se bazează algoritmul 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.

**Exemplul 7.4:** Pentru un exemplu foarte simplu cu cinci puncte în două dimensiuni să privim Fig. 15. 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 clusterelor și devin centroidul lor pentru moment.

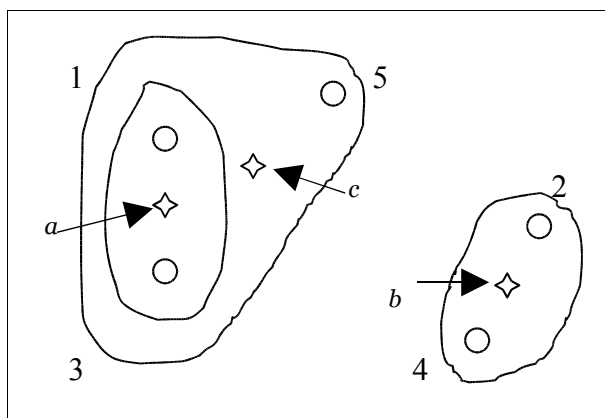


Figure 14: Un exemplu pentru algoritmul  $k$ -means

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  $c$ .

- 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 punctele la centroidul lor.

- 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 7.5. ilustrează acest lucru.

**Exemplul 7.5:** Să considerăm datele din Fig. 16. 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.

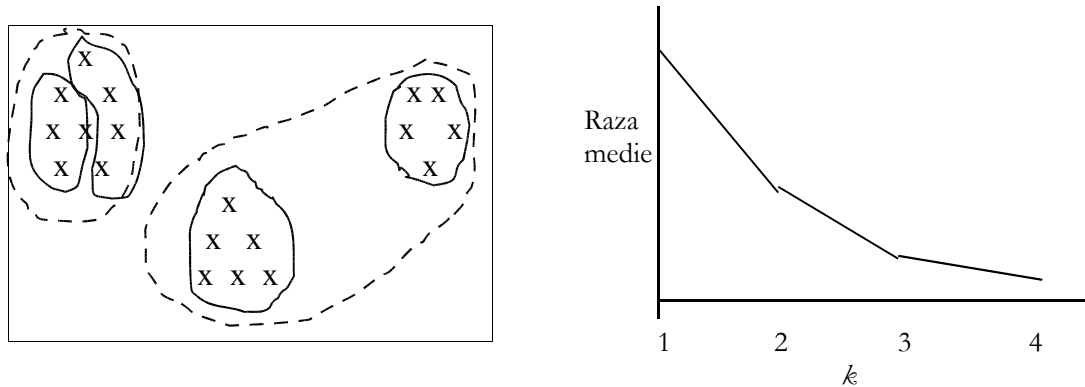


Figura 16. Descoperim că numărul corect de clustere este  $k=3$

Presupunem că apoi încercăm  $k=2$ . Unul dintre cele trei clustere 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.

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 Fig. 16. Totuși, dacă mărim  $k$  la 4 unul dintre adevăratele clustere va fi partiționat artificial în două clustere apropiate, așa cum arată liniile continue. 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.

## 7.6 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. Figura 17 arată cum ar trebui să arate datele pentru un cluster tipic în două dimensiuni. Un centroid, marcat cu +, are puncte împrăștiat de jur împrejur cu o deviație standard  $\sigma$  pe dimensiunea orizontală de două ori mai mare decât este aceasta pe dimensiunea verticală. Cam 70% din puncte vor fi în elipsa  $1\sigma$ ; 95% vor fi în cea  $2\sigma$ , 99.9% în  $3\sigma$  și 99.9999% în  $4\sigma$

## 7.7 Reprezentarea clusterelor în BFR

Cum a lucrat la Skycat, Usama Fayyad ("F"-ul din BFR) și-a reprezentat clusterelor ca pe niște galaxii, așa cum arată Fig. 18. Un cluster constă din:

1. Un nucleu central numit *Mulțimea de aruncat* (Discard set - DS). Mulțimea acestor puncte este considerată ca aparținând în mod sigur clusterului. Toate punctele din această mulțime sunt înlocuite de niște statistici simple, descrise mai jos. Notă: deși numite puncte "aruncate" 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.
2. Galaxii înconjurătoare, numite colectiv *Mulțimea comprimată* (Compression set - CS). 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.

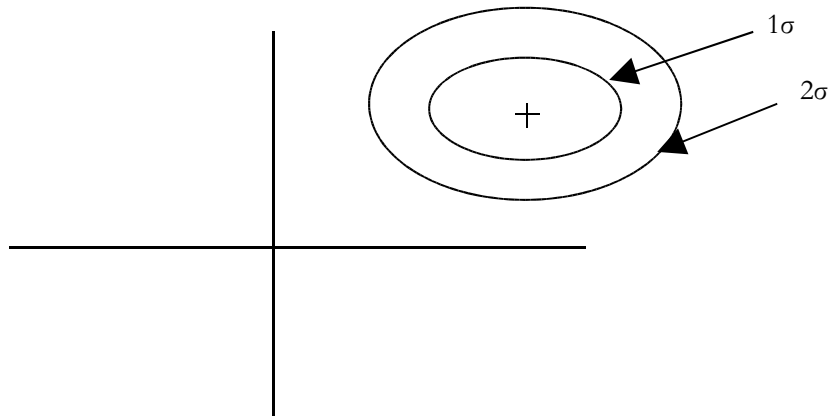


Figura 17. Se consideră că punctele dintr-un cluster au fost generate luând un centroid pentru cluster și adăugând la el în fiecare dimensiune o variabilă aleatoare normal distribuită cu media egală cu 0.

3. 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.

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  $SUMSQ_i$

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 și sunt mai convenabil de menținut pe măsură ce punctele sunt adăugate la cluster decât, să spunem, media și varianța în fiecare dimensiune. De exemplu:

- 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.

## 7.8 Procesarea unei memorii centrale de puncte în BFR

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. BFR propune două căi pentru a determina dacă un punct este suficient de aproape de un centroid pentru a fi luat în DS:

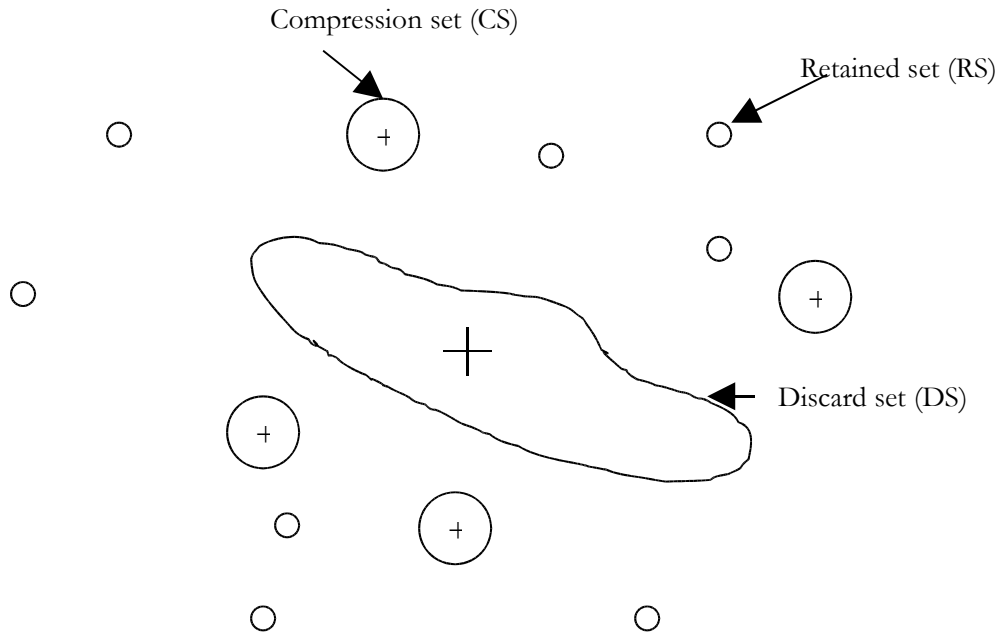


Figure 18: Exemple de clustere (DS), subclustere (CS) și puncte individuale (RS)

- (a) Se iau toate punctele pentru care *raza Mahalanobius* este sub un anumit prag, să zicem de patru ori deviația standard a clusterului. Raza Mahalanobius este în mare distanța până la centroid scalată în fiecare dimensiune cu  $\sigma_i$ , deviația standard în acea dimensiune. Mai precis, dacă  $\mu_i$  este media din dimensiunea  $i$ , atunci raza punctului  $\mathbf{y} = [y_1, y_2, \dots]$  este

$$\sqrt{\sum_i \left( \frac{y_i - \mu_i}{\sigma_i} \right)^2}$$

- (b) Pe baza numărului de puncte din diferite clustere, ne întrebăm dacă e probabil (să zicem la nivel de 95%) ca cel mai apropiat centroid curent să se mute în viitor suficient de departe de punctul  $\mathbf{y}$  și un alt centroid suficient de aproape de  $\mathbf{y}$  încât ultimul este atunci mai aproape de  $\mathbf{y}$  decât celălalt. Dacă nu este probabil ca cel mai apropiat centroid de  $\mathbf{y}$  să se schimbe atunci asignează  $\mathbf{y}$  la DS și plasează-l în clusterul celui mai apropiat centroid.
2. Se ajustează statisticile  $N$ ,  $SUM$  și  $SUMSQ$  pentru fiecare cluster pentru a include punctele tip DS adăugate.
  3. 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.
  4. 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.
  5. 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.