

Tehnici de multiindexare

Prof. dr. Ion SMEUREANU, asist. Marian DÂRDALĂ,
Catedra de Informatică Economică, A.S.E., Bucureşti

Articolul prezintă aspecte privind implementarea structurii de grid file pentru indexarea datelor multidimensionale.

Cuvinte cheie: grid file, partiție grilă, director grilă, scală liniară

Multiindexarea este una din metodele de organizare a datelor care permite regăsirea informațiilor după mai multe chei. În funcție de implementare, multiindexarea tratează diferențiat sau nediferențiat caracteristicile alese drept chei de indexare. Ne vom ocupa în continuare de o tehnică de multiindexare-grid file, care tratează nediferențiat cheile de selecție. Ea se bazează pe o reprezentare intuitivă a spațiului n -dimensional al caracteristicilor, căreia î se asociază o hartă de biți, matrice n -dimensională, semnalând prezența sau absența înregistrărilor care îndeplinesc cumulativ condițiile de selecție. Dimensiunea bitmap-ului va fi dată, deci, de numărul de atribute de identificare a înregistrărilor.

Presupunem un exemplu în care se dorește identificarea înregistrărilor după trei atribute. Se va construi un masiv tridimensional, în concordanță cu spațiul atributelor (S_1, S_2, S_3).

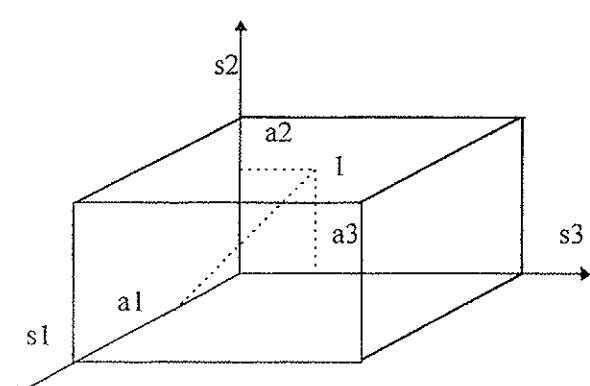


Fig. 1 Bitmap într-un spațiu tridimensional

Dacă se dorește găsirea înregistrării cu atributele (a_1, a_2, a_3) se consultă bitul corespunzător masivului tridimensio-

nal $B[a_1][a_2][a_3]$ și în caz că are valoare 1, înregistrarea există. Generalizând, se pot face operații de căutare a unei înregistrări, ceea ce se reduce la accesul direct în masiv a elementului $[a_1][a_2][a_3] \dots [a_k]$. Inserarea, respectiv ștergerea presupune poziționarea bitului pe valoarea 1, respectiv 0. Această metodă presupune ocuparea unui spațiu foarte mare de memorie, datorită varietății mari a valorilor de indexare, dezavantaj ce poate fi ameliorat aplicând diverse tehnici de compactare a masivelor.

Alt neajuns constă în faptul că în acest mod se verifică dacă există sau nu înregistrări care satisfac niște atribute, dar nu se precizează exact poziția lor în fișier. Pentru a elibera dezavantajul, se propune într-o primă etapă partiționarea spațiului căutării, prin alegerea unei diviziuni a domeniului. Alegerea diferențelor partiției depinde de distribuirea datelor, ce se presupune a nu fi întotdeauna uniformă, în timp ce atributele aferente partiției sunt independente.

O partiționare, în cazul tridimensional, rezultă prin fixarea unor intervale pe fiecare axă (X, Y, Z) ce divizează spațiul înregistrărilor sub forma:

$$U = (U_0, U_1, \dots, U_k), V = (V_0, V_1, \dots, V_m) \text{ și } W = (W_0, W_1, \dots, W_n).$$

În timpul operațiilor de lucru cu fișiere, partiții trebuie să fie modificate în concordanță cu inserările și ștergerile. Partiția $P = U * V * W$ este modificată prin alterarea unei componente la un moment de timp, datorită descompunerii unui interval în două, sau unirii a două intervale adiacente.

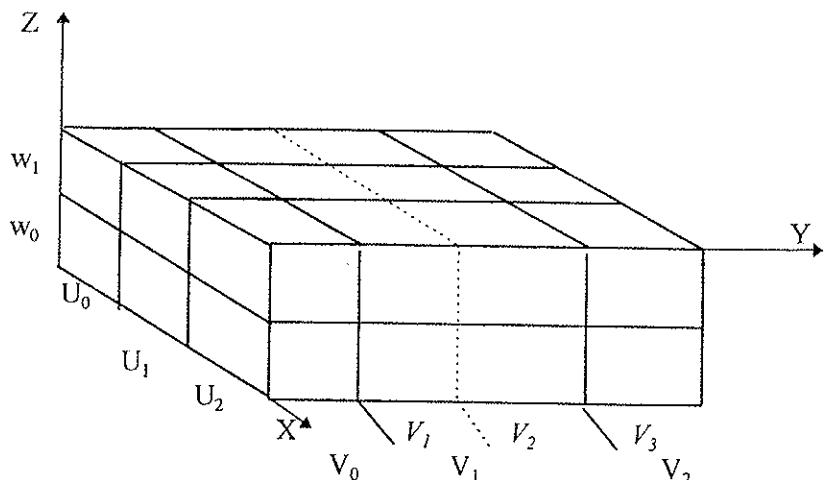


Fig. 2 Partițiile grilă

Metoda urmărește proiectarea unei structuri de fișier pentru stocarea și regăsirea înregistrărilor după mai multe atribută; structura se numește *fișier grilă*. Unitatea de stocare a înregistrărilor este *bucket*-ul și conține mai multe înregistrări. Pentru aceste fișiere, problema se reduce la definirea unei corespondențe între blocurile rezultante prin partiționare (blocuri grilă) și *bucket*-uri. Această asignare o realizează un *director grilă*.

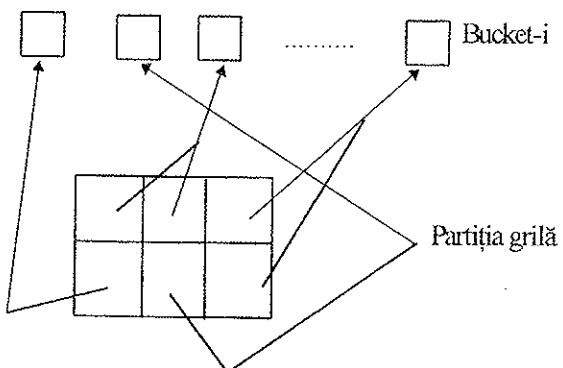
O posibilă structură de fișier grilă presupune:

- pentru o cerere punctuală să se folosească, în principal, două accese la disc;
- procesarea rapidă a cererii pe domenii mari, liniar ordonate;
- desfacerea și unirea blocurilor să implice doar două blocuri;
- menținerea unui nivel rezonabil de umplere medie a *bucket*-ului.

Un rol important revine directorului grilă care este de fapt un sistem de management a *bucket*-ilor, ceea ce implică:

- definirea unei clase de legături a blocurilor grilă cu *bucket*-ii;
- alegerea unei structuri de date pentru director, care să conțină asignarea curentă;
- găsirea unui algoritm eficient pentru a actualiza acea structură în mod direct, la o schimbare a asignărilor.

În figura 2 este un exemplu de asignare convexă a blocurilor la *bucket*-i.

Fig. 3 Legătura între directorul grilă și *bucket*-i

Fiecare partiție punctează un *bucket*. Există și situații când mai multe partiții punctează pe același *bucket*.

Un director grilă conține două părți:

- a) un masiv k -dimensional, care se numește masiv grilă, elementele lui fiind pointeri la *bucket*-i și sunt într-o corespondență unu la unu cu blocurile din partiție.

- b) k vectori, câte unul pentru fiecare atribut de indexare; un astfel de vector se numește *scală liniară*.

Pentru o mai bună înțelegere se va considera $k=2$, adică un spațiu al înregistrărilor, de forma $S = X * Y$. Un *director grilă* pentru un spațiu bidimensional se caracterizează prin:

- întregi $nx > 0$, $ny > 0$ ce precizează numărul de partiții a axei X , respectiv Y ;
- întregi $0 \leq cx < nx$, $0 \leq cy < ny$ ce precizează elementul curent din director și din blocul grilă curent;

- masivul bidimensional G ($0.., nx-1, 0.., ny-1$) numit masivul grilă;
- două masive unidimensionale X ($0.., nx$) și Y ($0.., ny$), ce reprezintă câte o scală liniară aferentă fiecărei dimensiuni a spațiului.

Operațiile ce trebuie să se definească pe structura directorului grilă sunt:

- accesul direct $G(cx, cy)$;
- accesarea pozițiilor vecine:

următorul / predecesorul x

$$cx = (cx \pm 1) \text{ mod } nx$$

următorul/predecesorul y

$$cy = (cy \pm 1) \text{ mod } ny$$

- jonctiunea, care se poate face atât între două partii de pe axa x cât și pe axa y . Astfel px , $1 \leq px < nx$, se unește cu următoarea parte ceea ce presupune renumirea tuturor elementelor care se află după px , cât și ajustarea scalei x care va avea un element mai puțin.
- descompunerea se poate face ca și operația precedentă pe scara x cât și pe

y . Având px , $0 \leq px \leq nx$, se va crea prin divizare un nou element $px+1$ și se renumește toate celulele ce succed lui px , scara x va avea un nou element.

Asignarea convexă a blocului grilă la *bucket* este exprimată prin următoarea relație: dacă $G(i', j')=G(i'', j'')$ atunci pentru toate valorile i, j cu proprietatea că $i' \leq i \leq i''$ și $j' \leq j \leq j''$ avem:

$$G(i', j')=G(i, j)=G(i'', j'').$$

În continuare se va prezenta un exemplu de folosire a directorului grilă în accesarea unei înregistrări în funcție de valorile a două atrbute de indexare: *an* cu valorile 0-2000 și *initială* cu domeniul a-z.

Se presupune că distribuția înregistrărilor a cauzat următoarea partiționare grilă:

$$x = (0, 1000, 1500, 1750, 1875, 2000)$$

$$y = (a, f, k, p, t)$$

Având cele două scale putem căuta înregistrarea cu atrbutele [1980,w].

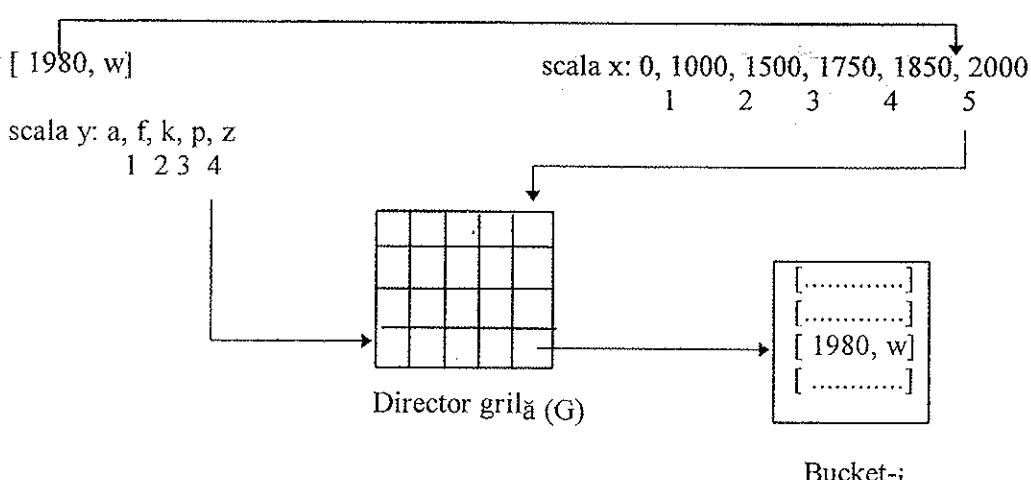


Fig. 4 Folosirea directorului grilă la identificarea *bucket*-ului pornind de la valorile cheilor

Determinarea *bucket*-ului care conține înregistrările cu valorile precizate pentru atrbutele de indexare se face în următorii pași:

- atrbutul 1980 este convertit în valoarea 5 din index printr-o căutare în scara x ;
- atrbutul w este convertit în valoarea 4 din index printr-o căutare în scara y ;
- indicii 4 și 5 permit accesul direct la elementul concret al gridului director

$G(4,5)$ de unde se extrage adresa *bucket*-ului ce conține înregistrarea dorită. Crearea fișierului grilă presupune un proces dinamic de construire, sincronizat cu cea a directorului grilă. Pentru aceasta se presupune inițial un singur *bucket* cu capacitatea de a reține trei înregistrări. Directorul grilă va conține, deci, un singur pointer la acest *bucket* unic.

Când *bucket*-ul *A* este plin (figura 5 a) și se încearcă introducerea unei noi înregistrări, atunci spațiul înregistrări-

lor este divizat și un nou *bucket* este disponibilizat (figura 5 b).

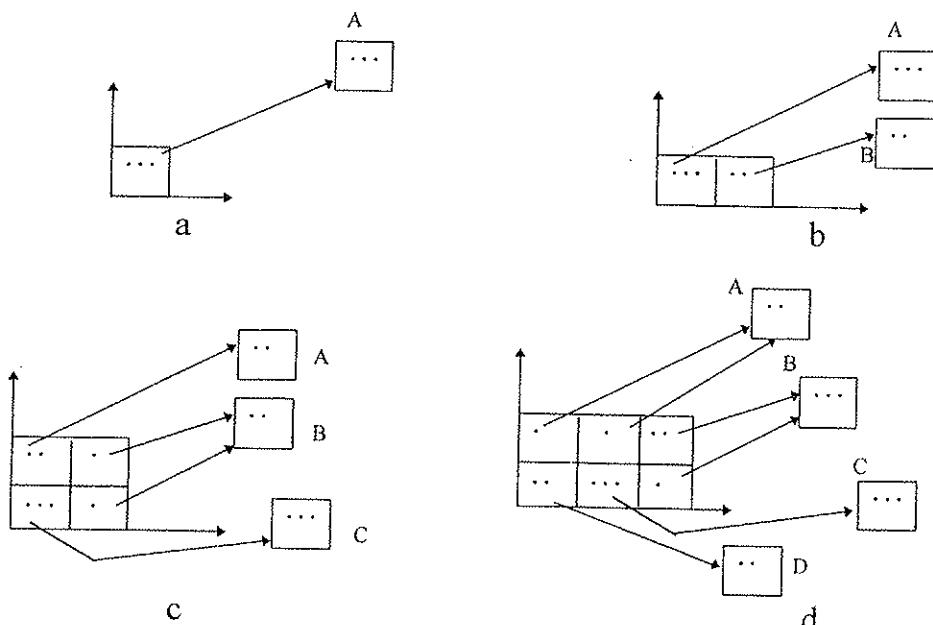


Fig. 5 Etapele de creare a fișierului grilă

Inregistrările sunt mutate din *bucket* vechi în cel nou, în funcție de apartenența la noua partiție ce se crează în directorul grilă. Dacă *bucket*ul *A* este din nou depășit, blocul grilă (jumătatea stângă a spațiului) este divizat după aceeași regulă (figura 5 c). Se presupune că se încarcă înregistrări și că *bucket*ul *C*, existent, va fi incapabil de a mai stoca înregistrări. Astfel se va crea un nou *bucket* (*D*) ca în figura 5 d.

În implementarea acestui tip de fișier se urmăresc anumite probleme ce sunt legate de mediu și în mare măsură de modul de realizare a directorului grilă. El poate fi implementat ca masiv *k*-dimensional sau folosind liste înălțuite.

Folosirea unui masiv are avantajul că accesarea elementelor se face mai ușor și se economisește memorie. Dezavantajul constă în greutatea realizării operației de inserare și divizare.

Listele înălțuite prezintă avantajul posibilității inserării, respectiv extragerii elementului din poziții arbitrale și deci a reconfigurării dinamice a structurii directorului grilă, dar presupun un

spațiu de memorie mai mare, necesar stocării pointerilor.

Pentru minimizarea timpului de răspuns la cereri, ultimile versiuni ale sistemelor de gestiune a bazelor de date (vezi Oracle Express Objects) au implementat și conceptul de date multidimensionale. Structura de tip grid file oferă una dintre implementările cele mai performante.

Bibliografie

1. Freeston, M., The BANG file: A new kind of grid file, Proceedings of the ACM SIGMOD Conference on Management of Data, San Francisco, May, 1987.
2. Nievergelt, J., Hinterberger, H., Sevcik, K.C., The Grid File: An Adaptable, Symmetric Multikey File Structures, Readings in DATABASE SYSTEMS, Morgan Kaufmann Publishers, San Francisco, California, 1994.
3. Salzberg, B., Grid file concurrency, Inf. Syst. 11,3,1986.