

DESCOPERIREA TIPARELOR FRECVENTE FĂRĂ GENERARE DE CANDIDAȚI

Cornelia GYÖRÖDI,

Robert GYÖRÖDI

Departmentul de Calculatoare, Facultatea de Electronică și Informatică,
Universitatea Oradea, România.

Tel: +40 (0) 259 432830 int.: 226 Fax: +40 (0) 259 432789

E-mail: rgyorodi@rdsor.ro

Rezumat: Informația implicită din bazele de date, în principal relațiile interesante de asociere dintre seturi de obiecte, care conduc la formarea de reguli de asociere, ar putea scoate în evidență tipare folositoare în activitățile de suport al deciziilor, previziuni financiare, politici de marketing, chiar și efectuarea de diagnostice medicale precum și în multe alte aplicații. Acest lucru a atras multă atenție în cercetările recente din domeniul data mining [1]. Așa cum este arătat în [2], descoperirea regulilor de asociere ar putea necesita scanări iterative ale unor baze mari de date, lucru care este costisitor din punct de vedere al timpului de procesare. Mulți autori și-au concentrat cercetările asupra problemei descoperirii eficiente a regulilor de asociere din bazele de date [2], [3]. Un algoritm de descoperire a regulilor de asociere, foarte influent, A priori [2], acest algoritm a fost dezvoltat pentru descoperirea regulilor de asociere, din baze de date tranzacționale de dimensiuni mari, bazat pe generarea de candidați și eliminarea celor care nu îndeplineau criteriile de suport și încredere impuse. Un pas important în îmbunătățirea performanțelor acestor algoritmi a fost făcut prin introducerea unei structuri de date compacte, numită *arbore de tipare frecvente* sau FP-tree [3], și a algoritmilor de descoperire asociați, FP-growth [3], care se baza pe descoperirea tiparelor frecvente fără generarea unor candidați suplimentari. Metoda introdusă de autor în [6] și [7], numită DynFPGrowth, îmbunătățește performanțele algoritmilor bazați pe FP-tree, prin reducerea numărului de scanări de la două la una singură și prin introducerea unui proces de actualizare dinamică a structurii FP-tree. Lucrarea prezintă, de asemenea, folosirea cadrului de dezvoltare și de comparare a algoritmilor de descoperire a regulilor de asociere, introdus în [8], în vederea validării metodei introduse și pe baza datelor obținute în urma unui studiu de performanță, care arată avantajele metodei introduse.

Cuvinte cheie: tipare frecvente, reguli de asociere, algoritmi, A priori, FP-Growth, DynFPGrowth

1. Introducere

Descoperirea regulilor de asociere are ca scop aflarea unui set de atribute comune, care aparțin unui număr mare de obiecte dintr-o bază de date. Descoperirea regulilor frecvente de asociere dintr-o bază de date, de dimensiuni mari, este o problemă complexă deoarece spațiul de căutare crește exponențial cu numărul de atribute din baza de date și cu obiectele bazei de date. Extragerea tiparelor frecvente joacă un rol esențial în descoperirea regulilor de asociere și a multor altor activități importante din domeniul minării datelor.

Cele mai multe din studiile anterioare, prezentate în [2], [4], [5] au adoptat o abordare bazată pe algoritmul A priori, care este întemeiat pe o euristică A priori antimonotonă [2]: *Dacă oricare tipar de lungime k nu este frecvent în baza de date, atunci supertiparul său de lungime $(k+1)$ nu poate fi frecvent niciodată.* Ideea esențială este de a genera, în mod iterativ, setul de tipare candidate de lungime $(k+1)$ din setul de tipare frecvente de lungime k (pentru $k \geq 1$) și de a verifica frecvența de apariție a acestora în baza de date.

Euristica A priori atinge performanțe bune prin reducerea dimensiunii seturilor de candidați. Cu toate acestea, în situații cu multe tipare frecvente, tipare lungi sau praguri de suport minim foarte joase, un algoritm bazat pe A priori poate, totuși, suferi de două probleme destul de costisitoare:

- este costisitoare manipularea unui număr mare de seturi de candidați; de exemplu, dacă există 10^4 seturi de articole frecvente de lungime 1, (1-articolset-uri), algoritmul A priori va trebui să genereze mai mult de 10^7 candidați de lungime 2 și să acumuleze și să verifice frecvențele de apariție ale acestora;
- este anevoioasă scanarea repetată a bazei de date și verificarea unui număr mare de candidați prin potrivire de tipare, în special pentru minarea tiparelor lungi.

Performanța metodelor gen A priori ar putea crește dacă s-ar putea evita generarea unui număr foarte mare de seturi de candidați. Această problemă a fost abordată de către Han în [3] prin trei aspecte.

În primul rând, este construită o structură de date nouă, denumită *arbore de tipare frecvente - frequent pattern tree* sau FP-tree, care este o structură de arbore-prefix extinsă, conținând informații cruciale, cantitative despre tipare frecvente. În arbore, vor exista noduri doar pentru articole frecvente de lungime 1, nodurile fiind aranjate astfel încât tiparele cu apariții mai frecvente vor avea o probabilitate mai ridicată în a partaja noduri decât cele mai puțin frecvente.

În al doilea rând, este dezvoltată o metodă de creștere a unui fragment de tipar, bazat pe FP-tree, care pornește de la un tipar frecvent de lungime 1 (ca și un tipar sufix inițial), examinează doar baza sa condițională de tipare – conditional pattern base (o sub-bază de date care constă din setul de articole frecvente care apar împreună cu tiparul sufix), construiește FP-tree-ul său condițional și efectuează minare, în mod recursiv, pe acesta. Creșterea tiparului este atinsă prin concatenarea tiparului sufix cu noile tipare generate din FP-tree-ul

condițional. Deoarece setul de articole (articolset-ul) frecvent din oricare tranzacție este întotdeauna codat în calea corespunzătoare a arborilor FP, creșterea tiparelor asigură completitudinea rezultatului. În acest context, metoda nu este gen A priori, adică o generare și verificare restricționată, ci este doar o verificare restricționată. Operațiile majore ale minării sunt acumularea de frecvențe și ajustarea contoarelor pe căile prefix, care sunt, în general, mult mai economice decât operațiile de generare de candidați și cele de potrivire de tipare efectuate în majoritatea algoritmilor gen A priori.

În al treilea rând, tehnica de căutare folosită în minare este bazată pe partiționare, o metodă mai degrabă tip *divide și cucerește*, decât gen A priori, *generare de jos în sus a combinațiilor de seturi de articole frecvente*. Aceasta reduce, în mod dramatic, dimensiunea bazei de tipare condiționale, generate la fiecare nivel de căutare ca și, de altfel, dimensiunea arborelui condițional corespunzător. Mai mult, transformă problema găsirii unor tipare frecvente lungi, în căutarea unora mai scurte și concatenând sufixele. Folosește cele mai puțin frecvente articole ca sufixe, ceea ce oferă o selectivitate bună. Toate aceste tehnici contribuie la reducerea substanțială a costurilor de căutare.

În paragraful 2 din această lucrare, voi prezenta metoda de construcție a unui arbore a tiparelor frecvente și voi descrie metoda FP-Growth [3], bazată pe o structură de arbore FP-tree, iar în paragraful 3, voi descrie o nouă metodă de construire a arborelui tiparelor frecvente cu performanțe promițătoare, introdusă în [6] pe care am denumit-o Dynamic FP-tree. Evaluarea performanțelor celor două metode am realizat-o în paragraful 4, prin folosirea cardului de comparare a algoritmilor de descoperire a regulilor de asociere, introdus în [8]. Concluziile sunt prezentate în paragraful 5.

2. Construcția arborelui tiparelor frecvente (FP-tree)

Fie $I = \{i_1, i_2, \dots, i_m\}$ un set de articole, iar $D = \langle T_1, T_2, \dots, T_n \rangle$ o bază de date tranzacțională, unde T_i ($i \in [1..n]$) este o tranzacție care conține un set de articole din I . **Supportul** (sau frecvența de apariție) a unui tipar A , care este un set de articole, este numărul de tranzacții din D care îl conține pe A . A este un *tipar frecvent* dacă supportul lui A nu este mai mic decât un suport minim *minsup* predefinit.

Fiind dată o bază de date tranzacțională D și un suport minim, *minsup*, problema găsirii setului complet de tipare frecvente este denumită **problema minării tiparelor frecvente**.

Proiectarea unei structuri de date compactă pentru minarea eficientă a tiparelor frecvente Han, în [3], se bazează pe un exemplu.

Exemplu 1: Fie baza de date tranzacțională D , primele două coloane din 2 – 1 și fie *minsup*=3.

O structură compactă poate fi proiectată pe baza următoarelor observații:

1. deoarece doar articolele frecvente vor juca un rol în minarea tiparelor frecvente, este necesară efectuarea unei scanări a bazei de date D pentru a identifica setul de articole frecvente;
2. dacă stocăm setul de articole frecvente ale fiecărei tranzacții într-o structură compactă, am putea evita scanarea repetată a bazei de date D ;
3. în cazul în care mai multe tranzacții au în comun un set de articole frecvente, acestea ar putea fi reunite într-un singur set, cu numărul de apariții înregistrate într-o variabilă *count*; este ușor de verificat dacă două seturi sunt identice dacă articolele frecvente în toate tranzacțiile sunt sortate după o anumită ordine dată;
4. dacă două tranzacții au în comun un anumit prefix, corespunzător unei ordonări a seturilor frecvente, părțile comune pot fi reunite folosind o structură prefix atât timp cât numărul de apariții *count* este înregistrat în mod corect. Dacă articolele frecvente sunt sortate descendent după frecvența lor de apariție, există șanse mai mari ca să avem mai multe șiruri prefix în comun.

Tabel 2 - 1. Baza de date cu tranzacții

TID	Articole cumpărate	Articole frecvente (ordonate)
100	F, a, c, d, g, i, m, p	f, c, a, m, p
200	a, b, c, f, l, m, o	f, c, a, b, m
300	b, f, h, j, o	f, b
400	b, c, k, s, p	c, b, p
500	a, f, c, e, l, p, m, n	f, c, a, m, p

Cu aceste observații se poate construi un *FP-tree* în următorul mod [3].

În primul rând, derivăm lista articolelor frecvente printr-o scanare a bazei de date D , $\langle(f:4), (c:4), (a:3), (b:3), (m:3), (p:3)\rangle$ (numărul de după: indică suportul), în care articolele sunt ordonate în ordine descrescătoare a frecvenței lor de apariție. Această ordonare este importantă deoarece fiecare cale din arbore va urma această ordonare. Pentru o mai ușoară înțelegere a algoritmului, în Tabel 2 - 1, coloana din dreapta arată această ordonare pentru fiecare tranzacție.

În al doilea rând, creăm rădăcina arborelui etichetat „root”, după care scanăm baza de date a doua oară. Scanarea primei tranzacții duce la construcția primei căi prin arbore: $\langle(f:1), (c:1), (a:1), (m:1), (p:1)\rangle$. De notat faptul că articolele frecvente din tranzacție sunt ordonate după ordinea listei articolelor frecvente. Pentru cea de a doua tranzacție, deoarece lista corespunzătoare de articole frecvente (în ordinea respectivă) $\langle f, c, a, b, m \rangle$ partajează un prefix comun $\langle f, c, a \rangle$ cu calea deja existentă $\langle f, c, a, m, p \rangle$, contorul fiecărui nod de-a lungul prefixului este incrementat cu 1, un nou nod $(b:1)$ este creat și legat ca fiu al nodului $(a:2)$ și un alt nod nou $(m:1)$ este creat și legat ca fiu al nodului $(b:1)$. Pentru a treia tranzacție, deoarece lista sa de articole frecvente $\langle f, b \rangle$ partajează doar nodul $\langle f \rangle$ cu subarborele prefix f , contorul lui f este incrementat cu 1 și un nou nod $(b:1)$ este creat și legat cu nodul $(f:3)$ ca fiu. Scanarea celei de a patra tranzacții conduce la crearea celei de a doua ramuri $\langle(c:1), (b:1), (p:1)\rangle$. Pentru ultima tranzacție, deoarece lista sa de articole frecvente $\langle f, c, a, m, p \rangle$ este identică cu prima, calea este partajată și contorul fiecărui nod de-a lungul căii este incrementat cu 1.

Pentru a facilita traversarea arborelui, este construit un tabel antet al articolelor, în care fiecare intrare arată spre apariția acestuia în arbore printr-un pointer la un nod al arborelui (cap de listă). Nodurile cu același articol sunt înlănțuite prin astfel de pointeri. Arborele rezultat după scanarea tuturor tranzacțiilor este arătat în figura 2-1.

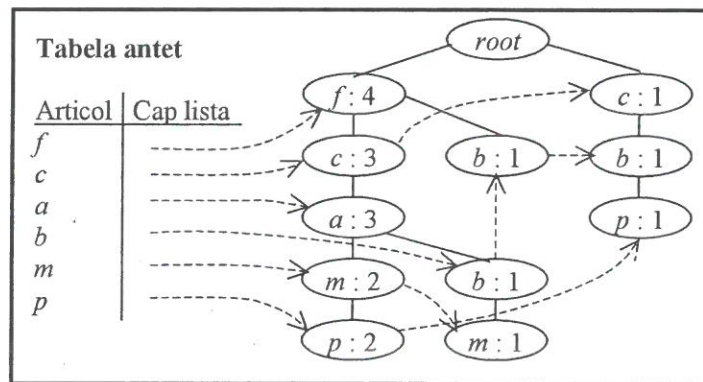


Figura 2 - 1. FP-tree pentru Exemplul 1

Acest exemplu conduce la următorul algoritm de construcție a arborelui tiparelor frecvente [3]:

Definiția 1 (FP-tree): Un arbore de tipare frecvente (FP-tree) este un arbore cu structura definită mai jos.

1. el constă dintr-un nod rădăcină, etichetat „root”, un set de subarbori prefix de articole ca fii ai rădăcinii și o tabelă antet pentru articolele frecvente;
2. fiecare nod din subarborele prefix de articole constă din trei câmpuri: *nume-articol*, *contor* și *înlănțuire-nod*, unde *nume-articol* se referă la articolul pe care acest nod îl reprezintă, *contor* înregistrează numărul de tranzacții reprezentate de porțiunea căii care ajunge la acel nod, iar *înlănțuire-nod* face legătura la următorul nod din FP-tree, care se referă la același articol sau conține NULL dacă nu mai există un alt astfel de nod;
3. fiecare intrare în *tabela antet* pentru articolele frecvente constă din două câmpuri (1) *nume-articol* și (2) *capul listei*, care arată spre primul nod din FP-tree care se referă la *nume-articol*.

Pe baza acestei definiții, avem următorul algoritm de construcție a arborelui FP-tree.

Algoritmul FP-tree (construcția FP-tree)

Intrare: O bază de date de tranzacții D și un suport minim s .

Ieșire: Arborele tiparelor frecvente (FP-tree) corespunzător.

Metoda: Arborele FP-tree este construit în următorii pași:

1. se scanează baza de date tranzacții D , se colectează setul de articole frecvente F și suportul corespunzător acestora, se sortează F în ordinea descrescătoare a suportului în L , lista articolelor frecvente;

2. se creează nodul rădăcină al unui arbore FP-tree, T și se etichetează cu „root”; pentru fiecare tranzacție $Trans$ din D se efectuează următoarele: se selectează și se ordonează articolele din $Trans$, conform ordonării lui L , adică în ordine descrescătoare a suportului; fie lista articolelor frecvente din $Trans$ $[p|P]$, unde p este primul element, iar P este restul listei; se apelează $insert_tree([p|P], T)$; funcția $insert_tree([p|P], T)$ este efectuată în felul următor: dacă T are un fiu N astfel încât $N.nume-articol = p.nume-articol$, atunci incrementează contorul lui N cu 1, altfel se creează un nou nod, cu contorul inițializat pe 1, părintele T și înlănțuire-nod înlănțuit cu nodurile având același $nume-articol$; dacă P nu este gol, se apelează în mod recursiv $insert_tree(P, T)$.

Din procesul de construcție a arborelui FP-tree putem observa că sunt necesare exact două scanări complete ale bazei de date D . Prima scanare colectează articolele frecvente, iar a doua scanare construiește arborele FP-tree. Costul inserării unei tranzacții $Trans$ în arbore este $O(|Trans|)$, unde $|Trans|$ este numărul de articole frecvente din $Trans$.

Dimensiunea unui arbore FP-tree este limitată de dimensiunea bazei de date corespunzătoare deoarece fiecare tranzacție va contribui cu, cel mult, o cale la arborele FP-tree, cu lungimea egală cu numărul de articole frecvente din acea tranzacție. Deoarece există, de multe ori, multe partajări ale articolelor frecvente între tranzacții, dimensiunea arborelui este, de cele mai multe ori, mult mai mică decât aceea a bazei de date originală. Spre deosebire de metodele gen A priori, care pot genera un număr exponențial de candidați în cel mai rău caz, în nici un caz, nu va fi generat un FP-tree cu un număr exponențial de noduri.

FP-tree este o structură foarte compactă, care stochează informația pentru minarea tiparelor frecvente. Deoarece o singură cale „ $a_1 \rightarrow a_2 \rightarrow \dots \rightarrow a_n$ ” în subarborele a_1 -prefix înregistrează toate tranzacțiile ale căror set maximal de articole frecvente este de forma „ $a_1 \rightarrow a_2 \rightarrow \dots \rightarrow a_k$ ” pentru orice $1 \leq k \leq n$, dimensiunea arborelui FP-tree este, substanțial, mai mică decât dimensiunea bazei de date și a dimensiunii setului de candidați generați pe parcursul minării regulilor de asociere.

Articolele din setul frecvent de articole sunt ordonate descrescător, în funcție de suportul lor. Articolele care apar mai frecvent sunt aranjate mai aproape de vârful arborelui FP-tree și, astfel, au o probabilitate mai mare de a fi partajate. Aceasta indică faptul că structura FP-tree este, în general, foarte compactă. De asemenea, experimentele arată că rezultă un arbore FP-tree destul de mic prin comprimarea unei baze de date destul de mari.

2.1. Descoperirea tiparelor frecvente folosind FP-tree. Metoda FP-Growth

În această secțiune, voi prezenta, conform **Error! Reference source not found.**, modul de explorare a informațiilor compacte, stocate în arbore, și voi descrie o metodă eficientă pentru minarea setului complet de tipare frecvente.

Algoritmul FP-Growth descoperă setul complet de articolset-uri frecvente, din baza de tranzacții D . Arborele FP-tree al bazei de date D conține informațiile complete ale bazei de date D în legătură cu minarea tiparelor frecvente, care au suportul peste pragul $minsup$.

Algoritmul FP-Growth pentru minarea tiparelor frecvente, utilizând arborele FP-tree, este descris în următoarea secvență:

Algoritmul FP-Growth (Minarea tiparelor frecvente cu FP-tree prin creșterea fragmentelor de tipare)

Intrare: Arborele FP-tree construit pe baza algoritmului FP-tree utilizând baza de date D și un suport minim s .

Ieșire: Setul complet de tipare frecvente.

Metoda: Apelul **FP-Growth**(FP-tree, null).

Procedure **FP-Growth**(Tree, α).

{

- (1). **if** Tree conține o singură cale P ;
- (2). **then for each** combinație (notată cu β) a nodurilor din calea P **do**;
- (3). generează tiparele $\alpha \cup \beta$ cu $support = suportul\ minim\ ale\ nodurilor\ din\ \beta$;
- (4). **else for each** a_i din tabela antet Tree **do** {
- (5). generează tiparele $\beta = a_i \cup \alpha$ cu $support = a_i.support$;

- (6). construiește baza condițională de tipare a lui β și
 apoi arborele FP-tree condițional a lui β notat $Tree_\beta$;
 - (7). **if** $Tree_\beta \neq \emptyset$
 - (8). **then call** **FP-Growth**($Tree_\beta$, β)
- }
- }

Dacă arborele conține o singură cale, tiparele generate sunt combinații ale nodurilor din cale, iar suportul fiind suportul minim al nodurilor din subcăi. Aceasta se realizează în liniile (1) - (3) ale procedurii FP-Growth. Altfel, se construiește baza condițională de tipare și se minează arborele FP-tree condițional pentru fiecare articolset frecvent a_i , iar suportul fragmentelor se ia ca fiind suportul articolset-urilor frecvente generate din baza condițională de tipare.

Procesul de minare bazat pe algoritmul FP-Growth scanează arborele FP-tree al bazei de date D o dată și generează o bază de tipare mică B_{a_i} pentru fiecare articol frecvent a_i , fiecare constând din seturi de căi prefix transformate ale lui a_i . Minarea tiparelor frecvente este apoi executată recursiv pe baza de tipare mică B_{a_i} cu construirea unui arbore FP-tree condițional pentru B_{a_i} . Așa cum a reieșit din analiza algoritmului FP-tree, un arbore FP-tree al unei baze de date este, în general, mult mai mic decât dimensiunea bazei de date. Similar, deoarece arborele FP-tree condițional, "FP-tree | a_i ", este construit pe baza de tipare B_{a_i} el ar trebui să fie mult mai mic și niciodată mai mare decât B_{a_i} . Mai mult, o bază de tipare B_{a_i} , este, în general, mult mai mică decât arborele său FP-tree original, deoarece ea constă din căile prefix transformate referitor la numai unul din articolele frecvente, a_i . Astfel, procesul de minare următor lucrează pe un set de baze de tipare și arbori FP-tree condiționali mult mai mici, în general. Aceasta este, în general, mult mai puțin costisitor decât generarea și testarea unui număr foarte mare de tipare candidat, fapt ce conduce la o eficiență mai mare a algoritmului FP-Growth.

Din algoritmul FP-Growth se observă că procesul de minare este un proces de tipul divide și cucerește (divide-and-conquer), și scala de micșorare este impresionantă. În general, factorul de micșorare este în jur de 20 ~ 100 pentru construirea unui arbore FP-tree dintr-o bază de date.

3. Construcția unui Dynamic FP-tree

Metoda introdusă în [6], numită Dynamic FP-tree, îmbunătățește performanțele algoritmilor bazați pe o structură de arbore, FP-tree, prin reducerea numărului de scanări ale bazei de date, de la două scanări la numai una singură și, de asemenea, permite un proces de actualizare dinamică a arborelui tiparelor frecvente, care se realizează pe o structură deja existentă.

Pentru a arăta modul de construcție a unui Dynamic FP-tree voi pleca de la construcția unui arbore FP-tree pe care am făcut câteva observații:

1. pentru o bază de date "logic" identică nu rezultă un arbore FP-tree unic; de exemplu, considerând baza de date cu tranzacții prezentată în Tabelul 2 - 1 și arborele corespunzător acesteia din figura 2-1, dacă primul articol din prima tranzacție ar fi c , și nu f , dar tranzacția ar conține aceleași articole, arborele FP-tree rezultat ar fi diferit;
2. procesul de construire a arborelui tiparelor frecvente are nevoie de două scanări ale bazei de date.

Pentru a depăși prima problemă, am propus o modificare a algoritmului de construcție a unui FP-tree, și anume aceea că, în pasul 1 al algoritmului de construcție al arborelui FP-tree, descris mai sus în secțiunea 2, setul de articole frecvente F ar trebui să fie sortat nu numai în ordine descrescătoare a frecvenței suport, ci și în ordine lexicografică (aceasta pentru articole care au același suport) rezultând L , lista articolelor frecvente. Aceasta ar trebui să asigure un arbore FP-tree unic pentru baze de date "logic" identice.

Pentru cea de-a doua problemă, am propus două abordări:

1. aplicarea directă a pasului 2, din algoritmul de construire a arborelui FP-tree, fără a scana prima dată baza de date pentru obținerea frecvențelor suport, utilizând numai o ordine lexicografică a articolelor într-o tranzacție; de asemenea, colectarea frecvențelor suport într-o tabelă header; suplimentar, avem pasul 3 în care tabela header este ordonată descrescător după suport și ordonată lexicografic și se reconstruiește arborele în conformitate cu noua ordine, dar numai prin scanarea arborelui;
2. proiectarea unui algoritm de reordonare dinamică a FP-tree astfel încât, aplicând direct pasul doi din algoritmul original, cu actualizarea mai întâi a frecvențelor suport în tabela header, păstrând ordinea

corespunzătoare atât a frecvențelor suport descrescătoare, cât și lexicografică a articolelor, iar în cazul în care există o "promovare" a cel puțin un articol, atunci să se efectueze o reordonare a arborelui pentru acele articole care au fost "promovate".

Ambele abordări vor produce un FP-tree care va reflecta toate tranzacțiile din baza de date. Această metodă de construire a arborelui și de stocare a arborelui este mai eficientă din punct de vedere a timpului consumat, în cazul unui număr mare de interogări cu diferite suporturi minime. Dar uneori arborele rezultat ar putea fi prea costisitor de memorat în întregime în memorie din punct de vedere a spațiului consumat de acesta.

Utilizând această metodă de construire dinamică a arborelui Dynamic FP-tree, nu trebuie să reconstruim arborele chiar dacă baza de date se actualizează prin adăugarea unor noi tranzacții. În acest caz, se va apela algoritmul, luând în considerare numai aceste noi tranzacții. Astfel, arborele se va modifica dinamic, prin adăugarea acestor noi tranzacții la forma lui anterioară.

Prin această metodă de construire a arborelui Dynamic FP-tree, se oferă un răspuns foarte rapid la orice interogare pe o bază de date, care se actualizează în mod continuu.

În continuare, voi descrie modificările care sunt făcute pe o structură FP-tree originală și, de asemenea, algoritmul de construcție a arborelui Dynamic FP-tree.

Deoarece în procesul de reordonare este nevoie să se modifice conținutul nodurilor din listă, pentru o mai mare ușurință, înlocuim lista simplă înlănțuită cu o listă dublu înlănțuită. Din acest motiv, introduc un pointer suplimentar *leg-ant* în fiecare nod din structura de nod, definită în secțiunea 2 (punctul 2 din Definiția 1), care conținea *nume-articol*, *contor* și *înlănțuire-nod*. Aceasta face legătura cu nodul anterior din arbore, care conține aceeași denumire de articol sau, dacă nu este nici unul, atunci va indica spre intrarea corespunzătoare din tabela header.

O altă modificare față de algoritmul original este aceea că headerul va avea două tabele, una corespunzătoare tabelii originale și a cea de-a doua, o tabelă *master* adițională, care constă dintr-un câmp, *nume-articol* și un câmp *contor* și va ține ordinea curentă de inserare în arborele FP-tree. Inserarea în tabela originală *original-ordered-table* poate fi făcută întotdeauna respectând ordinea impusă de suportul curent (ordine descrescătoare) și în ordine lexicografică. Orice articol nou va fi adăugat respectând această ordine.

Numim „promovare”, mutarea unui nod în arbore datorită creșterii frecvenței suport corespunzătoare. Orice „promovare” va fi detectată prin compararea tabelii *master* cu tabela originală, ordonată *original-ordered-table*. Pe această structură definesc o operație de *checkpoint*, care va sincroniza tabela *master* cu tabela originală, ordonată *original-ordered-table*, adică se copiază conținutul tabelii originale în tabela *master*. Aplicând operația de ordonare curentă, după inserarea unei tranzacții, vom obține un FP-tree actualizat. În conformitate cu această ordine curentă, arborele va fi reordonat dinamic.

Algoritmul de construcție a arborelui tiparelor frecvente Dynamic FP-tree este următorul:

Algoritmul Dynamic FP-tree (algoritmul de construire dinamică a arborelui tiparelor frecvente):

Intrare: O bază de date tranzacțională D și un suport minim s .

Ieșire: Arborele tiparelor frecvente (Dynamic FP-tree) corespunzător.

Metoda: Arborele tiparelor frecvente este construit în următorii pași:

1. Se creează nodul rădăcină T , al unui arbore FP-tree, și se etichetează cu „root”. Pentru fiecare tranzacție $Trans$ din baza de date, D execută următoarele:
 - a. adaugă articolele din tranzacția $Trans$ într-o tabelă header;
 - b. selectează și sortează articolele din tranzacția $Trans$ conform cu ordinea din tabela *master*; fie $[p | P]$ lista sortată de articole frecvente din $Trans$, unde p este primul element și P restul listei; apelează funcția $insert_tree([p | P], T)$;
 - c. funcția $insert_tree([p | P], T)$ este executată astfel: dacă T are un fiu N astfel că $N.item-name = p.item-name$ atunci incrementăm contorul lui N cu 1; altfel creăm un nod nou N și inițializăm contorul lui pe 1, legătura părinte a lui va indica spre T , iar legătura nodului va indica la nodurile cu aceeași denumire *item_name* prin structura de listă de noduri; dacă mai avem elemente în P , atunci se apelează recursiv funcția $insert_tree(P, N)$;
 - d. dacă se „promovează” un nod, atunci este nevoie de o reordonare și se apelează funcția $reorder(dynFPtree)$ pe arbore.

3.1. Reordonarea arborelui Dynamic FP-tree

Funcția de reordonare a arborelui, *reorder(dynFPtree)* se execută după cum urmează:

1. colectează toate articolele „promovate” (articolele care au frecvența suport modificată) într-o listă *reorderList*, ordonată după suport descrescător și în ordine lexicografică;
2. apelează *checkpoint ()* pentru a actualiza ordinea de inserare în FP-tree;
3. pentru fiecare articol din *reorderList* parcurge lista nodurilor și, pentru fiecare dintre aceste noduri, apelează funcția *moveUp(node)* pentru a pune acel nod în poziția corectă în FP-tree, în conformitate cu tabela *master*.

Funcția *moveUp(node)* se execută astfel:

1. repetă următoarea procedură până când nodul (*node*) și părintele său (*parent*) sunt în ordinea corectă:
 - a. preia părintele indicat de nodul părinte al nodului curent (*pparent*);
 - b. dacă *parent* are aceeași frecvență ca și *node*, elimină *parent* din lista *childNodes* a părintelui acestuia (*pparent*) și asignează-l lui *newNode*;
 - c. altfel efectuează următoarele acțiuni:
 - i. creează un nou nod *newNode* cu același articol ca și *parent*, dar având frecvența nodului *node*;
 - ii. inserează nodul *newNode* în lista nodurilor cu același articol;
 - iii. ajustează frecvența nodului *parent* prin scăderea frecvenței nodului *node*;
 - iv. elimină *node* din lista *childNodes* a lui *parent*.
 - d. înlocuiește lista fiilor din *newNode* cu lista fiilor lui *node* și actualizează legăturile părinte ale noilor fii cu noul părinte (*newNode*);
 - e. setează legătura părinte a lui *node* spre *pparent* (părintele original al părintelui nodului curent *node*), inițializează lista fiilor *childNodes* cu nodul *newNode* și setează legătura părinte a acestuia la *node*;
 - f. inserează nodul *node* în lista fiilor lui *pparent*.

3.2. Reducerea dimensiunii Dynamic FP-tree

Algoritmul prezentat asigură performanțe superioare în detrimentul consumului de memorie. Acest lucru se datorează faptului că, în funcția *moveUp* la pasul 1.f la un anumit nivel din FP-tree, am putea avea mai mult de un singur nod care să se refere la același articol. Acest lucru nu constituie nici o problemă pentru algoritmul de minare, dar duce la un arbore cu mai multe noduri decât necesar pentru reprezentarea compactă a tranzacțiilor din baza de date.

Din acest motiv, am dezvoltat un algoritm pentru combinarea a doi subarbori care au rădăcinile reprezentând același articol (*merge*).

Acest algoritm ar fi folosit în pasul 1.f din algoritmul *moveUp*, pasul în acest caz devenind:

- f. dacă avem, deja, un nod existent *existingNode* pentru articolul reprezentat de *node* în lista fiilor lui *pparent*, atunci apelează *merge(existingNode, node)* și continuă cu nodul *existingNode* ca și nod curent (*node*);
- g. altfel, inserează nodul *node* în lista fiilor lui *pparent*.

Funcția *merge(existingNode, node)* este implementată în următorul mod:

1. combină informația de frecvență prin incrementarea frecvenței suport a nodului existent *existingNode* cu frecvența nodului *node*;
2. elimină *node* din lista nodurilor cu același articol, dar fără modificarea legăturilor acestuia (acestea fiind necesare pentru continuarea algoritmului prin parcurgerea listei nodurilor corespunzătoare);
3. combină fiii nodului *node* cu fiii nodului existent *existingNode* prin executarea următoarelor acțiuni pentru fiecare fiu *child* al nodului:
 - a. dacă există un nod *sameItem* pentru același articol în nodul fiu al nodului *existingNode*, atunci combină cei doi prin apelarea recursivă a funcției *merge(sameItem, child)*;

- b. altfel, inserează nodul fiul *child* în lista fiilor nodului *existingNode* și setează părintele nodului fiu *child* la nodul *existingNode*.

4. Evaluarea experimentelor și studierea performanțelor

Luând în considerare rezultatele raportate în [8] și folosind cadrul de dezvoltare și comparare, am adaptat implementarea a doi algoritmi de descoperire a regulilor de asociere la acest cadru pentru a-l putea folosi și pentru a compara rezultatele obținute și pentru alți algoritmi, ca de exemplu A priori [2]. Primul algoritm implementat se baza pe o structură FP-tree, algoritm numit *FP-Growth* [3], iar celălalt bazat pe structura propusă Dynamic FP-tree, numit *DynFP-Growth*. Folosind acest cadru, am realizat câteva testări pe diverse baze de date. Pe algoritmul *DynFP-Growth* am realizat o optimizare prin folosirea funcției *merge*, descrisă în secțiunea 3.2, care are rolul de a reduce dimensiunea arborelui și am denumit algoritmul *DynFP-GrowthMerge*.

```
public class FPGrowth extends ARMiner {
    /** Arborele FP-tree */
    protected FPTree fpTree;
    /** Implementarea efectiva a algoritmului */
    public long mineFrequentItemsets() {
        ItemSet itemSet;
        ArrayList rawItemSet;

        frequentItemSets.clear();
        try {
            first();
            while((itemSet = getNextItemSet()) != null) {
                fpTree.addItemSet(itemSet);
            }
            fpTree.sortHeader();
            if (applyMinSupportToHeader)
                fpTree.hdr.applyMinSupport(getMinSuppCount());
            first();
            while((rawItemSet = getNextRawItemSet()) != null) {
                fpTree.addItemSet(rawItemSet);
            }
            fpGrowth(fpTree, frequentItemSets, new ArrayList(), getMinSuppCount());
        } catch (SQLException e) {
            e.printStackTrace();
        }
        frequentItemSets.internalSort();
        return 2;
    }
    protected void fpGrowth(FPTree fpTree, ItemSets frequentItemSets,
        ArrayList itemSet, long minSuppCount) {
        if (fpTree.singlePath()){
            /* Genereaza toate combinatiile din cale care satisfac cerintele */

```



```

}
else {
    /* Aplică recursiv algoritmul pe bazele condiționale de tipare */
}
}
/* ... */
}

```

Figura 4 - 1. Implementare FPGrowth

```

public class DynFPGrowth extends FPGrowth {
    /** Implementarea efectiva a algoritmului */
    public long mineFrequentItemsets() {
        ArrayList rawItemSet;
        long cnt;
        long crtSuppCnt;
        DynFPtree dynFPtree = (DynFPtree)fpTree;
        frequentItemSets.clear();
        try {
            cnt = 0;
            first();
            cnt++;
            while((rawItemSet = getNextRawItemSet()) != null) {
                crtSuppCnt = Math.round(cnt * getMinSupport());
                dynFPtree.addItem(new ItemSet(rawItemSet));
                dynFPtree.addItemSet(rawItemSet);
                if (!performReorderAtEnd && ((cnt % performReorderAtEvery) == 0)) {
                    if (dynFPtree.isReorderingNeeded(crtSuppCnt)) {
                        dynFPtree.reorder(crtSuppCnt);
                    }
                }
                cnt++;
            }
            if (performReorderAtEnd) {
                if (dynFPtree.isReorderingNeeded(getMinSuppCount())) {
                    dynFPtree.reorder(getMinSuppCount());
                }
            }
            // Nu trebuie sa sortam Header-ul deoarece este deja sortat !!!
            fpGrowth(dynFPtree, frequentItemSets, new ArrayList(), getMinSuppCount());
        } catch (SQLException e) {

```

```

    e.printStackTrace();
}
frequentItemSets.internalSort();
return l;
}
/* ... */
}

```

Figura 4 - 2. Implementare DynFPGrowth

Calculatorul pe care am realizat testele a fost un Pentium 4 la 1.7 GHz, cu 256 MBRAM sub sistemul Windows 2000. Pentru o evaluare cât mai realistă, am optat pentru baze de date, stocate pe servere SQL (Microsoft SQL 2000), accesul fiind realizat prin interfețele ODBC standard.

Pentru studierea performanțelor algoritmilor și pentru studierea scalabilității, am generat seturi de date, de la 10 000 până la 500 000 de tranzacții, și am folosit factori suport, începând de la 5% până la 40%. Orice tranzacție poate conține mai mult decât un articolset frecvent. Tranzacțiile pot avea un număr diferit de articole, de asemenea, și dimensiunea unui articolset frecvent diferă. Unele articolset-uri pot avea mai puține articole, altele pot avea mai multe. Având în vedere aceste observații, am generat seturi de date ținând cont de numărul de articole ale unei tranzacții, de numărul de articole dintr-un articolset frecvent etc. Parametrii necesari pentru a genera setul de date de test sunt definiți în Tabelul 4 - 1.

Tabel 4 - 1. Parametrii utilizați pentru generarea setului de date

D	Numărul de tranzacții
T	Media aritmetică a dimensiunii unei tranzacții
L	Numărul maxim de articolset-uri frecvente
N	Numărul de articole

Seturile de date de test sunt generate prin setarea numărului de articole la $N = 100$ și a numărului maxim de articolset-uri frecvente la $|L| = 3000$. De asemenea, se alege media aritmetică a dimensiunii unei tranzacții, $|T| = 10$.

Câteva rezultate comparative ale algoritmilor A priori, FP-Growth, DynFP-Growth, FP-GrowthMinSup și DynFP-GrowthMerge obținute pe diferite baze de date de test, raportate la un factor suport de 5%, sunt prezentate în Tabelul 4 - 2.

Tabel 4 - 2. Rezultate raportate la suport de 5%

Tranzacții (în mii)	Timp(secunde)				
	Minsuport 5%				
	A priori	DynFP-Growth	DynFP-GrowthMerge	FP-Growth	FP-GrowthMinSup
10	13.94	2.32	1.78	3.76	3.06
20	21.98	3.98	3.38	6.88	6.30
30	48.37	8.23	8.07	14.63	13.73
40	66.50	12.10	11.50	20.90	19.80
50	107.65	19.50	18.65	34.30	32.95
80	198.30	37.90	36.20	64.80	63.70
110	1471.40	55.00	53.30	95.50	93.90
150	3097.20	98.90	97.80	174.60	170.30
190	5320.60	152.70	150.00	273.60	268.00
300	9904.80	284.00	281.80	526.70	511.90
400	17259.20	458.10	449.80	849.70	849.20
520	20262.60	610.20	601.40	1150.70	1142.40

Rezultatele din Tabelul 4 – 2 arată că timpul de execuție al algoritmilor crește o dată cu dimensiunea bazei de date. Cea mai bună performanță o au algoritmi DynFP-GrowthMerge și DynFP-Growth. O performanță bună prezintă și algoritmul FP-Growth comparativ cu A priori, dar se observă că performanța lui DynFP-Growth este cu câteva ordine de mărime mai bună.

Din figura 4 – 3, se observă că timpul de execuție al algoritmilor FP-Growth și DynFP-Growth este constant raportat la aceeași bază de date, pentru un factor suport, care descrește de la 40% la 5%, pe când timpul de execuție al algoritmului A priori crește o dată cu descreșterea factorului suport. Pentru un factor suport mai mare de 30% pe o bază de date de 40 000 de tranzacții algoritmul A priori are performanțe asemănătoare cu DynFP-Growth și mai bune decât FP-Growth, dar pentru un suport mai mic de 20% performanța lui scade foarte mult. Astfel, timpul de execuție al algoritmului A priori ajunge să crească de până la 3 ori mai mult decât timpul de execuție al algoritmului FP-Growth și de până la 5 ori mai mult față de DynFP-Growth, pentru un suport de 5%.

Se observă că algoritmi FP-Growth și DynFP-Growth au o performanță bună din punct de vedere al timpului de execuție pentru diferite valori ale factorului suport, în timp ce algoritmul Apriori are o performanță mult mai slabă în special pentru factor suport mai mic de 20 %.

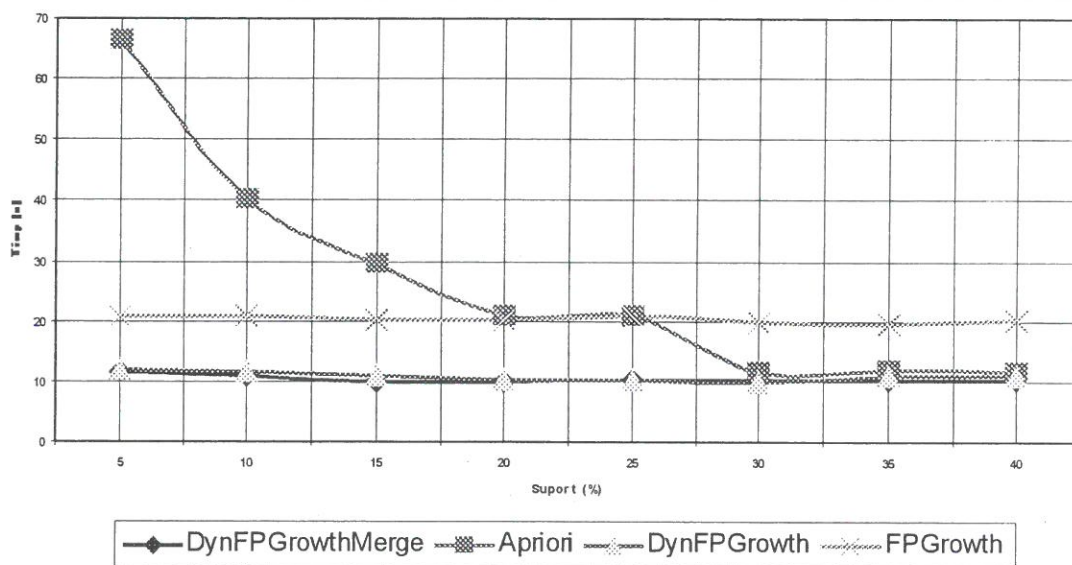


Figura 4 - 3. Scalabilitatea raportată la suport (D1 40K)

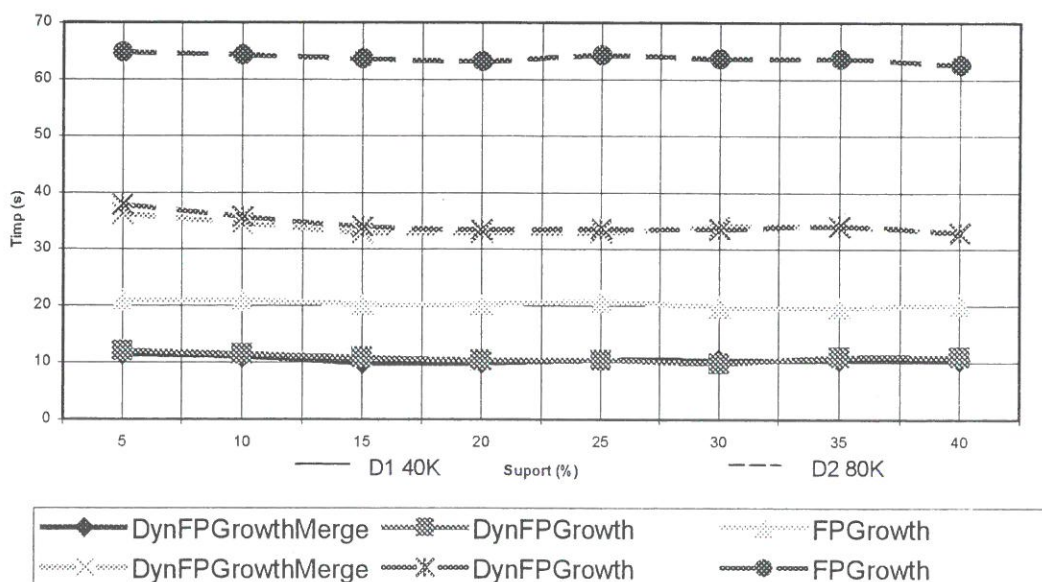


Figura 4 - 4. Scalabilitatea raportată la suport (D1 40K, D2 80K)

Scalabilitatea raportată la diferite valori ale suportului pe două baze de date (D1 = 40 000 și D2 = 80 000 de tranzacții) este prezentată în figura 4 - 4. Se observă că timpul de execuție al celor doi algoritmi depinde de dimensiunea bazei de date și nu de valoarea suport. Astfel, o dată cu creșterea dimensiunii bazei de date crește și timpul de execuție al algoritmilor, dar performanța algoritmului propus, DynFP-Growth este mai bună decât cea a algoritmului FP-Growth.

Din figura 4 - 5, se poate observa că atât algoritmul FP-Growth, cât și DynFP-Growth, au o performanță bună chiar și pentru un factor suport mic (minsup = 5%), și asta chiar pentru baze de date mari.

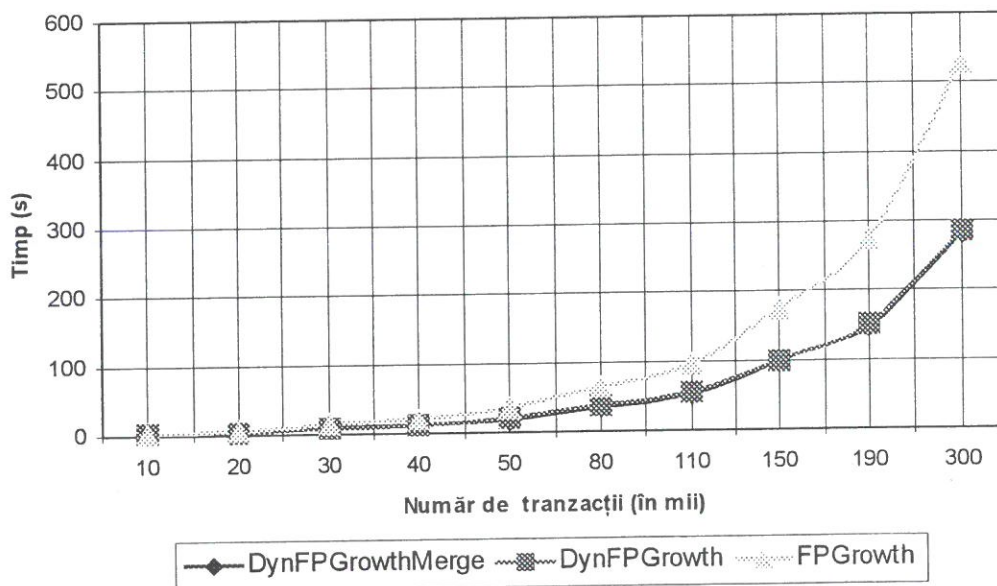


Figura 4 - 5. Timp de execuție pentru un factor suport de 5%

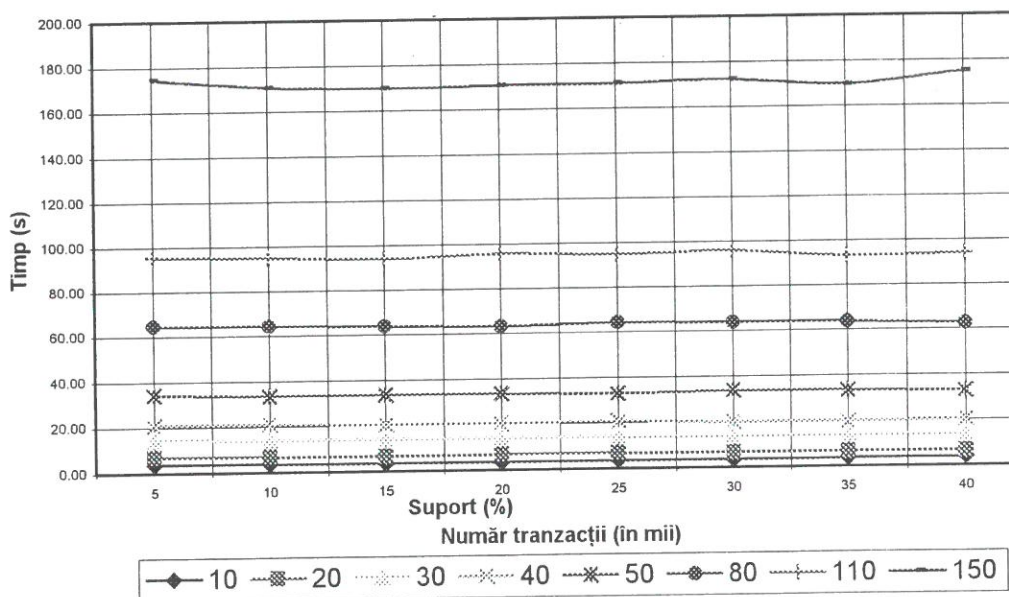


Figura 4 - 6. Scalabilitatea FP-Growth raportată la tranzacții/suport

Din figura 4 - 6, se observă că timpul de execuție al metodei FP-Growth nu depinde de suport, ci numai de dimensiunea bazei de date. În cazul metodei FP-GrowthMinSup, se ia în considerare suportul minim în construcția arborelui tiparelor frecvente, ceea ce face ca, pentru un factor suport mare, algoritmul FP-GrowthMinSup să fie mai performant, iar pentru un suport mic, de exemplu 5%, performanța lui să fie aproape identică cu al lui FP-Growth.

Metoda de construcție a arborelui articolset-urilor frecvente nu ia în considerare suportul, astfel arborele va conține toate tranzacțiile din baza de date și, în funcție de suport, se vor extrage din arbore articolset-urile frecvente,

care au frecvența mai mare decât suportul specificat. Această modalitate de construcție a arborelui permite aflarea articolset-urilor frecvente, pentru diferite suporturi, fără a mai necesita o nouă construcție a arborelui.

Prin aplicarea metodei propuse, de construire a arborelui tiparelor frecvente, DynFP-Growth, putem actualiza arborele prin adăugarea unor noi tranzacții (care sunt adăugate în baza de date ulterior, după ultima minare) la arborele deja construit, fără a reconstrui arborele din nou.

Numai o singură modificare trebuie făcută în algoritmul propus pentru o minare dinamică și anume, în pasul 1 al algoritmului Dynamic FP-tree:

1. încarcă arborele FP-tree, T , pentru fiecare tranzacție nouă $Trans$ din DB , execută pașii începând de la a , până la d , din algoritmul Dynamic FP-tree.

După modificarea arborelui tiparelor frecvente, putem aplica același algoritm FP-Growth pentru a efectua minarea articolset-urilor frecvente.

5. Concluzii

Din experimentele și studiile de performanță efectuate rezultă că algoritmul propus DynFP-Growth are performanțe mai bune decât algoritmul FP-Growth chiar și la baze de date mari. Timpii de execuție ai algoritmilor FP-Growth și DynFP-Growth nu depind de suport, ci numai de dimensiunea bazei de date: astfel, timpul crește o dată cu dimensiunea bazei de date, în timp ce performanța algoritmului A priori depinde atât de suport, cât și de dimensiunea bazei de date. Rezultatele minării, obținute prin metoda propusă DynFP-Growth, sunt compatibile cu cele obținute cu algoritmul FP-Growth pe o structură FP-tree originală **Error! Reference source not found.** Astfel, se obțin aceleași seturi de articole frecvente, dar ceea ce diferă este timpul de execuție a celor doi algoritmi, care este mai mic în cazul lui DynFP-Growth, acest lucru rezultând din experimentele prezentate în secțiunea 4.

Metoda propusă are performanțe mai bune deoarece se face numai o singură parcurgere a bazei de date, ceea ce minimizează timpul de creare a arborelui tiparelor frecvente. Această modalitate de construcție dinamică a arborelui Dynamic FP-tree poate fi folosită pentru a adăuga noi tranzacții la arborele deja existent, ceea ce metodele existente, bazate pe o structură de arbore nu o permit. Acesta este unul din avantajele principale, pe care o aduce noua metodă propusă, de construire dinamică a arborelui Dynamic FP-tree.

Abordarea propusă este mai potrivită în următoarele cazuri:

- când sunt baze de date mari, care trebuie interogate într-un timp scurt, și memoria principală nu este un factor determinat;
- când avem un arbore al tiparelor frecvente construit și ar trebui să-l actualizăm cu noi tranzacții fără să reconstruim întreg arborele tiparelor frecvente.

Bibliografie

1. **HAN, J., M. KAMBER:** Data Mining Concepts and Techniques, Morgan Kaufmann Publishers, San Francisco, USA, 2001, ISBN 1558604898.
2. **AGRAWAL R., R. SRIKANT:** Fast Algorithms for Mining Association Rules in Large databases. În Proc. of 20th Int'l conf. on VLDB, Santiago, Chile, September 1994, pp. 487-499.
3. **HAN, J., J. PEI, Y. YIN:** Mining Frequent Patterns without Candidate Generation. În: Proc. of ACM-SIGMOD, May 2000, ACM Press.
4. **PARK, J. S., M.S. CHEN, P.S. YU:** An Effective Hash-based Algorithm for Mining Association Rules. În: SIGMOD'95, pp. 175-186.
5. **SAVASERE, A., E. OMIECINSKI; S. NAVATHE:** An Efficient Algorithm for Mining Association Rules in Large Databases. În: VLDB'95, pp. 432-443.
6. **GYORODI, C.:** Contribuții la dezvoltarea sistemelor de descoperire a cunoștințelor. Teză de doctorat, Universitatea Tehnică "Politehnica" Timișoara, Romania, 2003.
7. **GYORODI, C., R. GYORODI, T. COFFEY, S. HOLBAN:** Mining Association Rules Using Dynamic FP-trees. În: Proc. of The Irish Signal and Systems Conference, University of Limerick, Limerick, Ireland, 30th June-2nd July, ISBN 0-9542973-1-8, pp. 76-82.
8. **GYORODI, R.:** A Comparative Study of Iterative Algorithms in Association Rules Mining". În: Studies in Informatics and Control, Volume 12 Number 3, Romania, 2003, ISSN 1220-1766, pp. 205-215.