

Studiu comparativ al algoritmilor de rezolvare a problemei Range Minimum Query

Anca Millio

Calculatoare si Tehologia Informatiei, anul II

1 Introducere

1.1 Descrierea problemei rezolvate

In domeniul IT, un Range Minimum Query (RMQ) rezolva problema gasirii valorii minime dintr-un subvector al unui vector de elemente comparabile. De exemplu: Fiind dat un vector A si doi indici, $i \leq j$, care este cel mai mic element dintre: $A[i]$, $A[i + 1]$, \dots , $A[j - 1]$, $A[j]$?

1.2 Exemple de aplicatii practice pentru problema aleasa

RMQ sunt utilizate in diferite cazuri, cum ar fi problema celui mai mic stramos comun (lowest common ancestor LCA) sau problema celui mai lung prefix comun (longest common prefix problem LCP).

1.3 Specificarea solutiilor alese

1. Sparse Table

Acest concept este utilizat pentru interogari rapide pe un set static de date (elementele nu se schimba). Este utilizata o functie care trebuie sa fie asociativa. Procesarea este facuta astfel incat interogariile primesc raspuns in mod eficient.

2. Sqrt-decomposition

Imparte intervalul dat in blocuri separate de cate \sqrt{n} elemente, apoi calculeaza minimul din fiecare bloc si stocheaza rezultatele. Este un algoritm care permite schimbarea valorilor din interval intre interogari. Raspunde la fiecare interogare in $O(\sqrt{N})$, realizeaza preprocesarea in $O(N)$.

3. Segment Tree

Frunzele arborelui sunt elementele vectorului dat, iar fiecare dintre nodurile interne reprezinta minimul tuturor frunzelor de sub acesta. De asemenea, acesta este un algoritm ce permite schimbarea valorilor vectorului in timpul interogarii. Raspunde la fiecare interogare in $O(\log N)$, realizeaza preprocesarea in $O(N)$.

1.4 Criteriile de evaluare pentru solutiile propuse

1. Se vor crea fisiere de input, care apoi vor fi pasate atat algoritmilor mentionati mai sus, cat si unui algoritm banal, apoi se vor compara fisierele de output, fiind evident ca output-ul algoritmului banal va fi corect, deoarece este usor de implementat si astfel sunt evitate greselile. Daca nu exista diferente, rezulta ca implementarile celor trei solutii sunt corecte. Datele vor fi generate aleator.
2. Pentru analiza eficientei si a complexitatii solutiilor, se vor crea fisiere de test cu dimensiuni variabile ale vectorului de numere, dar si a numarului de query-uri solicitate. Variatia duratei rularii in functie de marimea intrarii va fi stocata intr-un alt fisier, pe baza caruia vor fi apoi alcatuite grafice care sa ilustreze complexitatea fiecarui algoritm testat. Se vor alege vectori cu dimensiuni de 10, 100, 1000, etc elemente, iar pentru fiecare dintre aceste lungimi, cate un test cu 10 query-uri si unul cu 100 de query-uri. Numerele vor fi generate aleator.

2 Prezentarea solutiilor

2.1 Descrierea modului in care functioneaza algoritmii alesi

1. Sparse Table

Ideea pe care se bazeaza aceasta solutie este de a precalcula minimul tuturor subvectorilor de lungime 2^j , unde j variaza de la 0 la $\log n$, unde n este lungimea vectorului. Se realizeaza o matrice de cautare. Elementul $[i][j]$ va contine minimul intervalului incepand de la i si de marimea 2^j . De exemplu, cautarea $[0][3]$ va contine minimul din intervalul $[0, 7]$ (incepand cu 0 si dimensiunea 2^3).

Pentru orice interval arbitrar $[L, R]$, trebuie sa folosim intervale cu lungimi puteri ale lui 2. Se foloseste cea mai apropiata putere de 2. Trebuie intotdeauna sa facem cel mult o comparatie (compararea minimului a doua

intervale care sunt de lungimi puteri ale lui 2). Un interval incepe cu L si se termina cu "L + cea mai apropiata putere de 2". Celalalt interval se termina cu R si incepe cu "R - aceeaasi putere + 1". De exemplu, daca intervalul dat este [2, 10], comparam minimul a dou intervale [2, 9] si [3, 10].

2. Sqrt-decomposition

Se imparte intervalul [0, n-1] (ce reprezinta indicii elementelor din vectorul de lungime n), in blocuri de cate n elemente fiecare.

Se calculeaza minimul fiecarui bloc de dimensiune n si se stocheaza intr-un vector rezultatele.

Pentru a interoga un interval [L, R], luam minimul dintre valorile minime ale blocurilor care se afla in intregime in acest interval. Pentru blocurile care se suprapun partial cu intervalul dat, se parcurg liniar pentru a gsi minimul. La final se alege minimul dintre aceste valori.

3. Segment Tree

In Segment Tree, frunzele sunt elementele vectorului de intrare, iar fiecare nod intern reprezinta minimul tuturor frunzelor de sub el.

Se incepe cu un vector: v[0 . . n-1], care se imparte recursiv in jumatati. Se imparte de fiecare data segmentul curent in doua jumatati (daca acesta nu a devenit inca un segment cu lungimea 1), si apoi se apeleaza aceeaasi procedura pe ambele jumatati. Pentru fiecare astfel de segment, stocam valoarea minima intr-un nod al arborelui.

Odata ce arborele a fost construit, se pot face interogari. Urmatorul algoritim in pseudocod ilustreaza obtinerea rezultatului unei interogari.

$qs \rightarrow$ indicele de inceput al intervalului interogarii, $qe \rightarrow$ indicele de final al intervalului

```
int RMQ (nod, qs, qe)
{
    daca intervalul corespunzator nodului curent coincide cu
    intervalul interogarii(qs, qe)
        se returneaza valoarea nodului
    altfel daca intervalul corespunzator nodului curent este
    complet in afara interogarii(qs, qe)
        se returneaza infinit
```

```

    altfel
        se returneaza min(RMQ (copilul stng al nodului,
                             qs, qe), RMQ (copilul drept al nodului, qs, qe))
    }

```

2.2 Analiza complexitatii solutiilor

1. Sparse Table

Deoarece este nevoie de cel mult o comparatie pentru fiecare interogare, dupa cum am precizat mai explicit la sectiunea "Descrierea modului in care functioneaza algoritmi alesi", rezulta ca aceasta solutie are complexitatea interogarii de $O(1)$.

Complexitatea etapei de preprocesare a datelor poate fi calculata astfel:

In functia care creeaza matricea utilizata la interogare, bucla externa(for) are complexitatea $O(k)$, deoarece se alege un k minim, astfel incat $2^{k+1} > n$, asadar $O(k) = O(\log(n))$. Bucla interna(for) itereaza de la 0 la $n-1$ (unde n e lungimea vectorului) si contine doar o comparatie, deci are complexitatea $O(n)$. In concluzie, obtinem complexitatea $O(n * k) = O(n * \log(n))$.

Matricea formata are dimensiunea $n * \log(n)$, ceea ce determina o complexitate spatiala de $O(n \log(n))$.

2. Sqrt-decomposition

Aici vectorul este impartit in \sqrt{n} intervale de cate \sqrt{n} elemente ale vectorului, pri urmare complexitatea preprocesarii este de $O(\sqrt{n} * \sqrt{n}) = O(n)$, iar complexitatea spatiala este $O(\sqrt{n})$, deoarece este nevoie de crearea unui vector care stocheaza minimul din fiecare bloc.

Complexitatea interogarii este $O(\sqrt{n})$. Minimul blocurilor din interiorul intervalului dat este direct accesibil si pot exista cel mult \sqrt{n} blocuri de acest tip. Pot fi maxim doua blocuri care nu se suprapun complet cu intervalul dat si trebuie parcurse, asa ca va exista o complexitate de $2 * O(\sqrt{n})$ pentru acest proces. Asadar, complexitatea totala va fi $O(n)$.

3. Segment Tree

Complexitatea pentru constructia arborelui(preprocesarea) este $O(n)$. Numarul de noduri din arbore este de $2n-1$, iar valoarea fiecarui nod este cal-

culata o singura data in constructia arborelui.

Complexitatea interogarii este $O(\log n)$. Pentru a interoga minimul unui interval, se proceseaza cel mult doua noduri la fiecare nivel, iar numarul de nivele este $O(\log(n))$ (inaltimea arborelui este $\log(n)$).

Complexitatea spatiala este determinata de stocarea arborelui, deci apartine $O(n)$.

2.3 Prezentarea principalelor avantaje si dezavantaje pentru solutiile luate in considerare

1. Sparse Table

– Avantaje:

Utilizeaza o structura de date simpla, aceasta fiind matricea pentru aflarea raspusurilor la interogari.

Un alt avantaj consta in viteza de raspuns la interogari ($O(1)$), fiind cea mai rapida solutie dintre cele trei, din acest punct de vedere.

– Dezavantaje:

Daca vectorul de input sufera modificari, costul repetarii preprocesarii va reprezenta un dezavantaj, fiind cel mai ineficient ca durata dintre cei trei algoritmi comparati ($O(n \log(n))$).

2. Sqrt-decomposition

– Avantaje:

Implementarea este usoara, in comparatie cu celelalte doua solutii, avand nevoie de o structura de date foarte simpla, mai exact vectorul in care se stocheaza minimele blocurilor vectorului de input ($O(\sqrt{n})$).

Preprocesarea este mai putin costisitoare decat in cazul "Sparse table".

– Dezavantaje:

Complexitatea interogarii reprezinta un dezavantaj, fiind cea mai lenta comparativ cu celelalte solutii alese ($O(\sqrt{n})$).

3. Segment Tree

– Avantaje:

Interogările se fac mai rapid decât în cazul Sqrt-decomposition, iar memoria utilizată și viteza de preprocesare depășesc performanța Sparse Table, ceea ce duce la concluzia că acest algoritm este recomandat în cazul în care se fac modificări ale vectorului de intrare între interogări, dar se dorește și o viteză relativ bună de răspuns la interogări.

– Dezavantaje:

Cantitatea de cod este mai mare decât la celelalte soluții, făcând implementarea mai anevoioasă.

3 Evaluare

3.1 Descrierea modalității de construire a setului de teste folosite pentru validare

În afara de cele trei soluții, am implementat și un algoritm cu rezolvarea banală, pentru a putea compara corectitudinea celor trei soluții.

Pentru realizarea testelor am creat un fișier(generator.cpp) în care am implementat un algoritm care generează un vector de numere, dar și un vector de perechi de numere(pentru interogări) în mod aleator, apoi am scris într-un fișier bash(Tests), comenzi care dau ca parametrii codului generator, numărul de elemente din vector și numărul de interogări dorite și redirecționează vectorul de rezultate într-un fișier corespunzător din folder-ul "out", precum și durata interogărilor în alt fișier(*Algo_Time*). Acest ultim fișier este utilizat la construirea graficelor care ilustrează performanța comparativă a soluțiilor alese.

Am realizat teste pe diferite lungimi ale vectorului de intrare(10, 100, 1000, ...), precum și număr variabil de interogări(10, 100) pentru fiecare lungime a vectorului, apoi am aplicat fiecare test fiecărui algoritm(algo1 = Sparse Table, algo2 = Sqrt-decomposition, algo3 = Segment Tree, *basic_algo* = soluția banală).

3.2 Specificațiile sistemului de calcul pe care au fost rulate testele

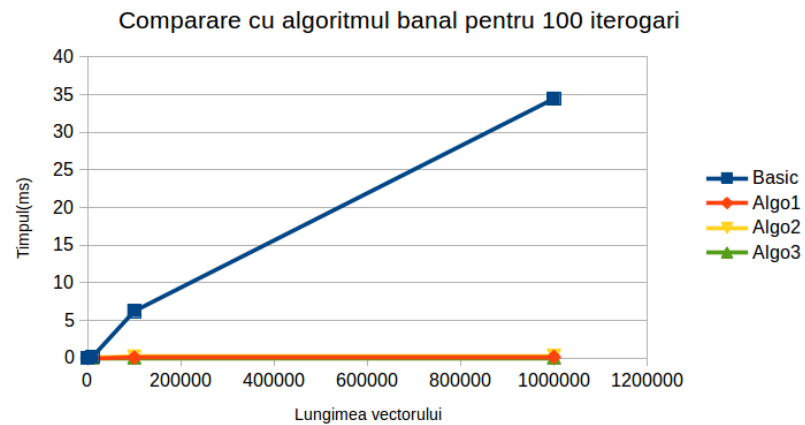
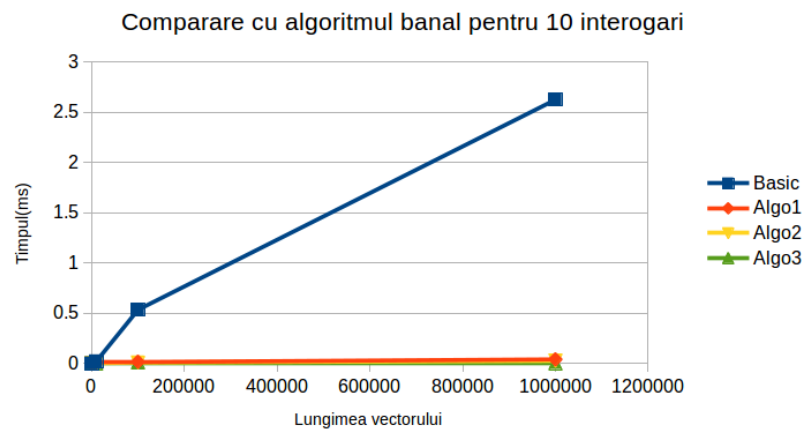
Procesor:

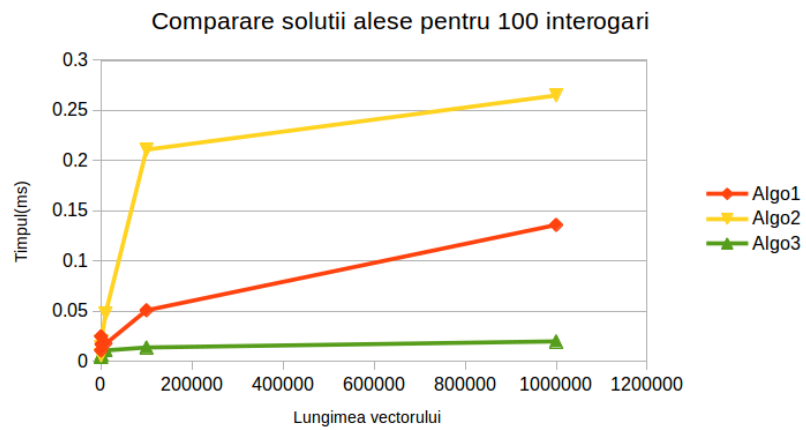
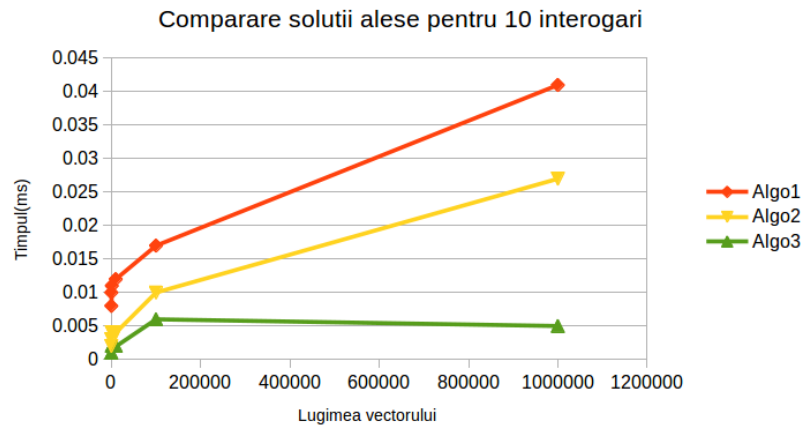
Intel(R) Core(TM) i7-7500U CPU @ 2.70GHz

Memorie disponibila:
8 GB

3.3 Ilustrarea, folosind grafice, a rezultatelor evaluarii solutiilor pe setul de teste

Se poate observa ca cele trei solutii sunt mult mai bune decat algoritmul banal pentru oricate interogari.





3.4 Prezentarea succinta a valorilor obtinute pe teste

Rezultatele sunt cele asteptate in privinta performantei algoritmilor 2 si 3, atat pentru 10, cat si pentru 100 de interogari, algoritmul 3 fiind mai rapid decat 2.

Este neasteptat faptul ca pentru 10 interogari, algoritmul 1 este cel mai ineficient, iar pentru 100 de interogari, acesta se situeaza pe locul doi. Ne-am fi asteptat ca el sa fie pe locul intai in toate situatiile.

Analizand codul, am observat ca pentru fiecare interogare se calculeaza un logarithm, iar acest calcul consuma foarte mult timp. Pentru a verifica aceasta ipoteza, am comentat toate liniile din functia query a algoritmului 1 (Sparse Table), in afara de calculul logarithmului si am observat ca timpii de executie sunt sensibili egali cu varianta in care codul este complet. Pentru a aduce performanta la rezultatele asteptate, ar trebui construit un vector care sa stocheze rezultatele calculelor logarithmilor respectivi in faza de preprocesare.

4 Concluzii

In practica nu exista o solutie perfecta(cea mai buna). Este esential ca aceasta sa fie aleasa in functie de caracteristicile datelor de intrare.

Daca vectorul de input se modifica des(de exemplu de 100 de ori pentru 1000 de interogari) sau se dau foarte putine interogari, atunci ar fi optim sa se aleaga solutia cu cel mai bun timp de preprocesare. Segment Tree si Sqrt-decomposition au cea mai buna complexitate de preprocesare(fata de Sparse Table) mai exact $O(n)$, deci preferabil ar fi sa alegem Segment Tree, deoarece acesta are o performanta mai buna la interogari.

Daca vectorul de input nu se modifica, dar se fac multe interogari(de la ordinul zecilor de mii in sus), cel mai potrivit ar fi Sparse Table, deoarece interogariile se fac cel mai rapid($O(1)$), comparativ cu celelalte doua solutii.

5 Bibliografie

Data ultimei accesari a urmatoarelor surse este: 13.12.2018.

1. <https://web.stanford.edu/class/cs166/lectures/01/Slides01.pdf>
2. <https://www.geeksforgeeks.org/range-minimum-query-for-static-array/>
3. <https://www.geeksforgeeks.org/segment-tree-set-1-range-minimum-query/>