

CAPITOLE SPECIALE DE INFORMATICĂ

Arborei kd

1 Noiembrie 2018

Un arbore kd (prescurtare de la „arbore k-dimensional”) este o structură de date arborescentă pentru o colecție de puncte dintr-un spațiu k-dimensional $S = \underbrace{A_0 \times \dots \times A_{k-1}}$ unde A_0, \dots, A_{k-1} sunt mulțimi total ordonate¹ de valori, care poate fi folosită pentru a rezolva problema următoare:

Se dă un punct $X \in S$ și o mulțime finită de puncte $M \subset S$

Să se găsească punctul $P \in M$ care este cel mai apropiat de punctul X .

Un arbore kd pentru M este un arbore binar care memorează câte un punct din M în fiecare nod, și satisfac condițiile următoare:

- Pentru fiecare nod X la nivelul n , dacă (a_0, \dots, a_{k-1}) este punctul reținut în nodul X și $i = n \bmod k$, atunci

- $a'_i \leq a_i$ pentru toate punctele (a_0, \dots, a_{k-1}) stocate în subarborele stâng al lui X ,
- $a_i < a'_i$ pentru toate punctele (a'_0, \dots, a'_{k-1}) stocate în subarborele drept al lui X .

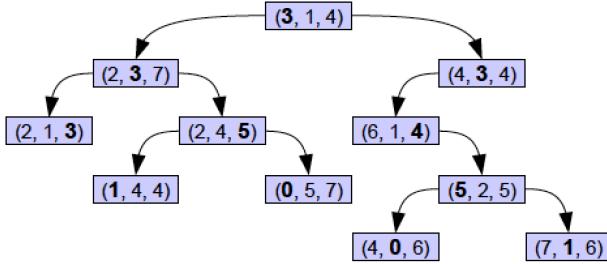
Deci, un arbore kd este o generalizare a unui arbore binar de căutare în care poziția cheii de căutare dintr-un punct depinde de nivelul nodului care reține punctul respectiv: cheia de căutare dintr-un punct X stocat într-un nod de la nivelul n este coordonata $(n \bmod k)$ a lui X ; presupunem că coordonatele punctelor sunt indexate pornind de la 0.

De exemplu, dacă $k = 3$, $A_0 = A_1 = A_2 = \mathbb{R}$ și

$$S = \{ (0, 5, 7), (1, 4, 4), (2, 1, 3), (2, 3, 7), (2, 4, 5), \\ (3, 1, 4), (4, 0, 6), (4, 3, 4), (5, 2, 5), (6, 1, 4), (7, 1, 6) \}$$

atunci un arbore kd care reține nodurile mulțimii S este

¹O mulțime A este total ordonată dacă există o relație binară \leq pe A care este **totală** (pentru toți $a, b \in A$, fie $a \leq b$ sau $b \leq a$), **reflexivă** ($a \leq a$ pentru toți $a \in A$), **antisimetrică** (pentru toți $a, b \in A$, dacă $a \leq b$ și $b \leq a$ atunci $a = b$) și **tranzitivă** (pentru toți $a, b, c \in A$, dacă $a \leq b$ și $b \leq c$ atunci $a \leq c$).



1 Construcția unui arbore kd balansat

Fie L o listă de puncte k -dimensionale și o axă $i \in \{0, 1, \dots, k - 1\}$. Putem defini recursiv arborele kd $T(L, i)$ pentru lista de puncte L și axa inițială i , în felul următor:

Cazul de bază: Dacă L este lista vidă atunci $T(L, i)$ este arborele vid.

Cazul recursiv: Fie $L' = [P_1, \dots, P_n]$ rezultatul sortării listei L în ordine creșcătoare a valorilor coordonatei i a punctelor, $m_0 = \lfloor n/2 \rfloor$, și

$$m_1 := \max\{j \mid m_0 \leq j \leq n \text{ și } P_{m_0}, P_j \text{ au aceeași valoare a coordonatei } i\}.$$

Atunci $T(L, i)$ re

- punctul P_{m_1} stocat în nodul rădăcină,
- subarborele stâng $T([P_1, \dots, P_{m_1-1}], (i+1) \bmod k)$, și
- subarborele drept $T([P_{m_1+1}, \dots, P_n], (i+1) \bmod k)$.

Arborele kd construit în acest fel este balansat dacă $m_0 = m_1$ în toți pașii recursivi.

2 Găsirea unui nod în un arbore kd

Presupunem că T este un arbore kd cu axa inițială i . Pseudocodul algoritmului de găsire a nodului din T care reține un punct $A = (a_0, \dots, a_{k-1})$ este:

```

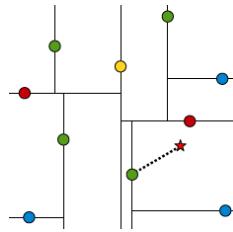
FINDPOINT( $T, i, A$ )
if  $T$  este gol then return 'nod negasit'
( $a'_0, \dots, a'_{k-1}$ ) := punctul din rădăcina lui  $T$ 
if  $a'_j = a_j$  for  $0 \leq j < k$  then return rădăcina lui  $T$ 
if  $a_i \leq a'_i$  then
    return FINDPOINT( $T.\text{left}, (i+1) \bmod k, A$ )
else /*  $a_i > a'_i$  */
    return FINDPOINT( $T.\text{right}, (i+1) \bmod k, A$ )

```

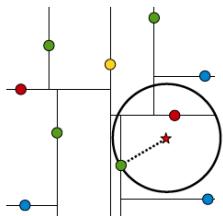
3 Găsirea celui mai apropiat vecin al unui punct

Presupunem că T este un arbore kd cu axa inițială i , și că vrem să găsim punctul stocat în T care este cel mai aproape de un punct test $A = (a_0, \dots, a_{k-1})$.

Intuiția din spatele algoritmului prezentat aici este următoarea. Să presupunem că avem o bănuială care este nodul din T cel mai apropiat de punctul A . De exemplu, să presupunem că punctul A (punctul de test) este cel indicat de o steluță, și că credem că vecinul lui cel mai apropiat este cel conectat la steluță cu linie întreruptă.



Remarcă: Dacă ar exista un punct mai apropiat de A decât punctul cu care este conectat prin linie întreruptă, atunci acel punct ar trebui să fie în cercul cu centrul în A care trece prin bănuiala curentă:



Observații

- În un spațiu k -dimensional am avea o hipersferă în loc de cerc, pe care am numi-o **hipersferă candidat**.
- Remarca precedentă ne permite să eliminăm subarborei ai arborelui kd din spațiul de căutare a celui mai apropiat vecin: Dacă toată hipersfera este într-o singură parte a unui hiperplan de separare, atunci cel mai apropiat vecin nu poate fi în partea cealaltă a hiperplanului de separare, și prin urmare putem omite căutarea unui vecin în subarborele pentru punctele din cealaltă parte a hiperplanului.

De pildă, în exemplul ilustrat cercul se află complet în dreapta liniei verticale care trece prin rădăcina arborelui. Prin urmare, orice punct stocat în subarborele stâng al rădăcinii nu poate fi mai aproape decât nodul cu care este conectat A prin linie întreruptă.

- Criteriul de detectie dacă o hipersferă cu centrul (a_0, \dots, a_{k-1}) și raza d nu este complet într-o parte a hiperplanului definit de ecuația $x_i = \text{curr}_i$ este foarte simplu: $|\text{curr}_i - a_i| < d$.
- Dacă coordonatele punctelor sunt valori reale, se obișnuiește să se presupună că distanța dintre două puncte $X = (x_0, \dots, x_{k-1})$ și $Y = (y_0, \dots, y_{k-1})$ este distanța euclidiană,

adică:

$$distance(X, Y) = \sqrt{(x_0 - y_0)^2 + \dots + (x_{k-1} - y_{k-1})^2}.$$

Pe baza acestor remarci putem defini algoritmul următor (pseudocod) de găsire a punctului din T care este cel mai aproape de A :

```

1NN( $T, i, A$ )
Maintain a global best estimate of the nearest neighbor, called guess
Maintain a global value of the distance to that neighbor, called bestDist
guess := Null
bestDist :=  $\infty$ 
1NNAUX( $T.\text{root}, i, A$ )
return guess
```

where

```

1NNAUX(curr,  $i, A$ )
if curr is empty then return
if  $distance(A, \text{curr}) < \text{bestDist}$  then
    guess := curr
    bestDist :=  $distance(A, \text{curr})$ 
/* recursive search of  $A$  in current kd-tree curr */
if  $a_i \leq \text{curr}_i$  then
    search := 'left'
    1NNAUX(curr.left, ( $i + 1$ ) mod  $k, A$ )
else
    search := 'right'
    1NNAUX(curr.right, ( $i + 1$ ) mod  $k, A$ )
/* Check if the candidate hypersphere crosses this splitting hyperplane */
if  $|\text{curr}_i - a_i| < \text{bestDist}$  then
    if search = 'left' then
        1NNAUX(curr.right, ( $i + 1$ ) mod  $k, A$ )
    else
        1NNAUX(curr.left, ( $i + 1$ ) mod  $k, A$ )
```

4 Găsirea celor mai apropiati K vecini ai unui punct

Algoritmul 1NN poate fi generalizat pentru a găsi cei mai apropiati K vecini ai unui punct test $A = (a_0, \dots, a_{k-1})$ într-un arbore kd T cu axa inițială i : în loc să ținem evidența unei singure bănuieri globale **guess**, putem reține intr-o structură globală cei mai apropiati K vecini descoperiți până la momentul respectiv. O structură de date potrivită pentru acest lucru este o coadă cu priorități și capacitate limitată (engl. *bounded priority queue*, BPQ), al cărei comportament este descris în Appendix.

Algoritmul generalizat se numește k NN (prescurtare de la „ k nearest neighbors”). Pseudocodul acestui algoritm este indicat în figura 4.

```

kNN( $T, i, A$ )
Maintain a BPQ of the candidate nearest neighbors, called bpq
Set the maximum size of bpq to  $K$ 
kNNAUX( $T.\text{root}, i, A$ )
return bpq

```

where

```

kNNAUX( $\text{curr}, i, A$ )
if curr is empty then return
/* add curr to bpq */
enqueue q into bpq with priority  $distance(A, \text{curr})$ 
/* recursively search the half of the tree that contains the test point  $A$  */
if  $a_i < \text{curr}_i$  then
    search := 'left'
    kNNAUX(curr.left,  $(i + 1) \bmod k, A$ )
else
    search := 'right'
    kNNAUX(curr.right,  $(i + 1) \bmod k, A$ )
 $p :=$  priority of max.-priority element of bpq
if bpq is not full OR  $|\text{curr}_i - a_i| < p$  then
    if search = 'left' then
        kNNAUX(curr.right,  $(i + 1) \bmod k, A$ )
    else
        kNNAUX(curr.left,  $(i + 1) \bmod k, A$ )

```

Figure 1: Algoritmul k NN.

Temă de laborator

1. (O structură de date pentru arbori kd) Folosiți C++ pentru a implementa `kdTree`, o structură de date pentru arbori kd care rețin o mulțime de puncte implementate ca instanțe ale clasei

```

struct Point {
    float x[3]; // coordonatele punctelor
}

```

și implementați metoda

- `makeKdTree(vector<Point> v)` care crează un arbore balansat kd din punctele reținute în un vector `v`

2. (Problema k NN) Scrieți un program care face următoarele:

- (a) Crează un arbore kd balansat T pentru o mulțime de n puncte, ale căror coordinate sunt citite din un fișier text `points.txt` cu conținutul următor:

Are n linii, și fiecare linie conține trei numere în virgulă flotantă separate prin spațiu, care reprezintă coordonatele unui punct.

De exemplu, `points.txt` ar putea avea conținutul următor:

```
1.0 4.0 5.14
-17.3 25 6.42
0 0 0.3
3.0 4.0 5.0
```

(b) Așteaptă ca utilizatorul să introducă de la consolă

- valoarea lui K (un număr întreg pozitiv) și
- coordonatele punctului de test A (trei numere `float`)

și returnează K noduri din T care sunt cele mai apropiate de A .

Utilizați structurile `Point` și `BPQ` disponibile de pe site-ul cursului, și structura de date `kdTree` implementată pentru prima temă de laborator.

A Cozi cu priorități și capacitate limitată

O coadă cu priorități și capacitate limitată (BPQ) este asemănătoare cu o coadă cu priorități obișnuită, cu excepția faptului că are o capacitate maximă care reprezintă numărul de elemente care pot fi reținute în ea. Ori de câte ori se adaugă un element nou la ea, dacă coada este plină, se elimină din ea elementul care are prioritatea cea mai mare. De exemplu, să presupunem că avem o BPQ cu capacitatea de 5 elemente, care reține elementele următoare:

Valoare	A	B	C	D	E
Prioritate	0.1	0.25	1.33	3.2	4.6

Ce se întâmplă dacă vrem să inserăm în ea elementul F cu prioritatea 0.4? Deoarece coada este plină, se va inseră elementul F după ce se elimină elementul cu cea mai mare prioritate (E). În final, conținutul cozii va fi

Valoare	A	B	F	C	D
Prioritate	0.1	0.25	0.4	1.33	3.2

Apoi, să presupunem că vrem să inserăm elementul G cu prioritatea 4.0 în această coadă. Deoarece prioritatea lui G este mai mare decât cea mai mare prioritate a unui element din coadă, coada de priorități rămâne neschimbată (nu se va întâmpla nimic).

Observație: o implementare a BPQ, concepută să rețină instanțe ale clasei `Point` care reprezintă puncte din \mathbb{R}^3 , este dată în `BPQ.cpp` and `BPQ.h`, și poate fi descărcată de pe site-ul acestui curs.