

On cherche à implémenter la structure de données abstraite Dictionnaire, qui représente un ensemble de valeurs, chacune étant associée à une clé dans un ensemble totalement ordonné. On s'intéresse en particulier aux opérations :

- Recherche: clé x dictionnaire \rightarrow valeur
- Insérer: valeur x dictionnaire \rightarrow dictionnaire
- Supprimer: clé x dictionnaire \rightarrow dictionnaire

ainsi qu'à des opérations de recherche d'une valeur en fonction du rang de sa clé.

\rightarrow Cette structure de données est d'une grande importance pour le stockage d'informations, et l'optimisation des opérations ci-dessus a donc beaucoup de conséquences pratiques : gestion d'une base de données, stockage d'informations sur un disque dur...

\rightarrow On suppose dans la suite que les clés sont distinctes.

I] Recherche dans une liste ou un tableau.

1] Recherche par clés :

\rightarrow Dans une liste non-triée :

On recherche une clé en parcourant la liste : complexité $O(n)$ au pire et en moyenne; de même pour la suppression; insertion en $O(1)$.

\rightarrow Efficace pour peu de données

\rightarrow Dans un tableau trié :

Recherche dichotomique: Si T est trié, on peut faire une méthode diviser-pour-régner :

- Comparer la clé cherchée avec celle de $T[\lfloor \frac{|T|}{2} \rfloor]$
- Si elle est égale, renvoyer $T[\lfloor \frac{|T|}{2} \rfloor]$
- Si elle est inférieure (resp. supérieure), rechercher l'élément dans $T[1.. \lfloor \frac{|T|}{2} \rfloor]$ (resp. $T[\lfloor \frac{|T|}{2} \rfloor + 1.. |T|]$)

Recherche en temps $O(\log n)$, Insérer et Supprimer nécessitent $O(\log n)$ opérations pour déterminer l'emplacement, puis le coût d'ajout ou de suppression d'un élément, soit $O(n)$.

\rightarrow Utile pour les tableaux statiques (Rechercher uniquement)

Prop 1: (Optimalité) Les algorithmes de recherche par comparaison pour les tableaux triés ont une complexité en temps $\Omega(\log_2 n)$ (en moyenne et au pire).

\rightarrow Optimisation de Insérer / Supprimer ?

Si on a beaucoup de telles opérations à faire, la complexité $O(n)$ peut être un problème. Nous verrons en II et III comment améliorer ces performances.

2] Recherche par rangs: [Cor] chap 9

\rightarrow Recherche du min, du max, de la médiane, du k^o plus grand

Faire ce type de recherche dans un tableau trié est constant. Pour une liste non triée, on a donc un algorithme en $O(n \log n)$ pour faire ces recherches. Ce n'est pas optimal.

• Rechercher le min (ou le max):

Comparer un minimum courant à toutes les clés.

\rightarrow temps $O(n)$ optimal

• Rechercher le k^o plus grand à l'aide d'un diviser pour régner

Prendre une clé quelconque, séparer le tableau T en

$$T_+ = \{b \in T, b > a\} \quad T_- = \{b \in T, b < a\}$$

Si $|T_-| < k$, $|T_+| > |T| - k$, appeler récursivement sur T_-

Si $|T_-| > k$, $|T_+| < |T| - k$, appeler récursivement sur T_+

sinon, renvoyer a

\rightarrow complexité $O(n^2)$ au pire cas, $O(n)$ en moyenne.

\rightarrow Recherche du plus petit élément plus grand qu'une clé.

On peut construire un algorithme similaire à la recherche dichotomique. min en $O(n)$.

Les algorithmes de recherche par rangs sont facilement implémentables avec des arbres (II), et efficace.

[Fon] Partie 3
9.2.1

9.3.1

Hyp: des ordonnées

9.3.4

9.1

9.2

[?]

II] Arbres binaires de recherche [Cor] chap 12

1) Définition et implémentation [Froi] chap 10
[Sed] chap 14

Def 2: Un arbre binaire de recherche (ABR) est un arbre binaire étiqueté tel que pour tout nœud v , les éléments du sous-arbre gauche (resp. droit) sont inférieurs (resp. supérieurs) à v : (cf Annexe 1)

On représente un dictionnaire par un ABR en étiquetant les sommets par les clés:

Recherche x A: $v = \text{racine}(A)$

Si $x = v$, terminé.	} Complexité $O(\text{haut}(A))$
Si $x < v$, Recherche x gauche(A)	
Si $x > v$, Recherche x droite(A)	

Insérer x A:

Si $A = \emptyset$, renvoyer x	} Complexité $O(\text{haut}(A))$
sinon, $v = \text{racine}(A)$	
Si $x < v$, Insérer x gauche(A) Si $x > v$, Insérer x droite(A)	

Supprimer x A:

Rechercher x A	} Compl. $O(\text{haut}(A))$
Si x est une feuille, supprimer x .	
Si x n'a qu'un fils y , remplacer x par y Si x a deux fils g et d , remplacer x par le plus grand élément du sous-arbre g	

Il est également aisé de trouver le Successeur ou le Prédécesseur d'un élément, de trier, ou de trouver le $k^{\text{ème}}$ élément dans un ABR.

→ **Pire cas:** Arbre dégénéré. Opérations en $O(n)$.

Se produit si on construit l'arbre à partir d'un tableau trié.

→ **Meilleur cas:** Arbre équilibré: Opérations en $O(\log_2 n)$.

Prop 3: Les opérations Recherche, Insérer, Supprimer se font en moyenne en temps $O(\log_2 n)$ dans un ABR.

Si le pire cas est inévitable, on peut chercher à maintenir l'arbre équilibré.

2] Équilibrage et AVL [Froi] Partie 3, chap 11

→ Rotations

Pour équilibrer un ABR, on va utiliser des opérations de rotation (représentées en Annexe 2) sur les arbres: une réassignation d'un nœud et d'un de ses fils (et éventuellement petit-fils) et de leurs sous-arbres afin de modifier les hauteurs relatives des sous-arbres \rightarrow se fait en temps constant.

→ AVL

Def 4: Un AVL est un ABR tel que pour tout nœud, la différence de hauteur entre ses sous-arbres gauche et droit est d'au plus 1.

Pour implémenter un AVL, on stocke dans un nœud les hauteurs de ses sous-arbres gauche et droit.

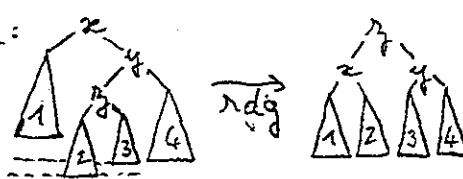
Prop 5: Si A est un AVL à n nœuds, alors $\text{haut}(A) = O(\log_2(n))$

La recherche, insertion et suppression fonctionnent comme pour un ABR, mais pour insérer et supprimer, il faut remonter dans l'arbre à partir de l'élément modifié et rééquilibrer le sous-arbre dans lequel on se trouve:

Cas 1:



Cas 2:



Prop 6: Les opérations Recherche, Insérer, Supprimer se font au pire en temps $O(\log_2 n)$ dans un AVL.

Remarque 7: Il existe d'autres types d'arbres permettant d'améliorer les ABR: arbres rouge-noir, B-arbres (non binaire).

3] Applications

→ ABR optimaux et compression [Cor] 15.5, [Knu]

Si p_i est la probabilité que l'élément i soit recherché, et q_j ($0 \leq j < n$) celle qu'un élément $j < x < j+1$ soit recherché,

Notation
N(A) = nombre de nœuds
H(A) = hauteur

[Froi] 10.4.2

nombre d'éléments

on cherche à construire un arbre minimisant le temps moyen d'une recherche \rightarrow ceci peut se faire en $O(n^3)$.

\rightarrow Dans le cas où les p_i sont nuls, ceci permet de réaliser le codage de Huffman pour compresser un texte.

[PS] \rightarrow Intersections de segments dans le plan: DEV

7.2
(p. 28)

Étant donné N segments dans le plan, on cherche à lister leurs points d'intersection: Utilisation d'ABR pour représenter la structure du problème \rightarrow si $I = \text{nb d'intersections}$: $O((I+N) \log_2 N)$.

III Tables de Hachage [Cor; Ch. 11]

Dans cette partie, on suppose les clés dans $[1, m]$.

1) Un tableau pour les recherches toutes

Adressage direct

On déclare un tableau de taille m , chaque valeur sera envoyée dans la case de sa clé.

La recherche, l'insertion, la suppression se font alors en $O(1)$.

\rightarrow si le nombre de valeurs est petit devant m , on utilise inutilement beaucoup trop de mémoire...

Hachage

On choisit $a \in [1, m]$ et $h: [1, m] \rightarrow [1, a]$ f. de hachage. Une valeur de clé k sera envoyée sur l'adresse $h(k)$.

Si les valeurs sont envoyées sur des cases différentes, les coûts de recherche, insertion et suppression sont toujours $O(1)$.

\rightarrow si $a < m$, il y a forcément des clés k, l telles que $h(k) = h(l)$.

def 8: On dit qu'il y a collision si deux valeurs doivent être envoyées sur la même case.

Gérer les collisions

• adressage: si une adresse est déjà remplie à l'insertion, on trouve une adresse vide algorithmiquement.

recherche/suppression: si $T(h(\text{clé})) \neq \text{valeur}$, on suit l'algorithme d'insertion.

• chaînage: (utilisé par la suite): utiliser des listes (cf Annexe 3)

2) Faire un bon hachage

But: minimiser/empêcher les collisions

def 9: Un hachage $h: [1, m] \rightarrow [1, a]$ est dit uniforme si pour chaque valeur de clé tirée, et pour tout $i \in [1, a]$, $P(h(k) = i) = \frac{1}{a}$

\rightarrow impossible en pratique sans hypothèses sur les distributions des clés; on peut cependant bien s'en approcher.

Avec un tel hachage, le coût moyen d'une opération insertion/suppression/recherche pour n valeurs présentes est $O(\frac{n}{a})$

\rightarrow Si la fonction est connue d'avance, un utilisateur malveillant peut tenter de saturer une adresse (ex: attaque d'un site web).

Pour régler ce problème, on prendra à la création du tableau une fonction au hasard parmi une classe universelle:

def 10: Une classe universelle \mathcal{H} est un ensemble de fonctions de hachage $[1, m] \rightarrow [1, a]$ tel que pour chaque paire de clés (k, l) , $|\{h \in \mathcal{H}, h(k) = h(l)\}| \leq \frac{|\mathcal{H}|}{a}$

ex: soit p premier, $p > a$; on pose

$$\mathcal{H}_{p,a} = \{h \mapsto ((xk + y) \bmod p) \bmod a \mid x, y \in \mathbb{Z}/p\mathbb{Z}, x \neq 0\}$$

Hachage parfait. [Cor; 11.5]

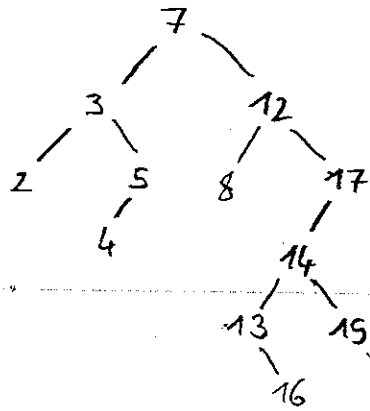
Il est possible, en utilisant un double hachage, de ranger n valeurs dans une structure de taille mémoire $O(n)$ telle que la recherche dans cette structure soit en $O(1)$. La création de cette structure se fait en $O(n)$ DEV

Conclusion

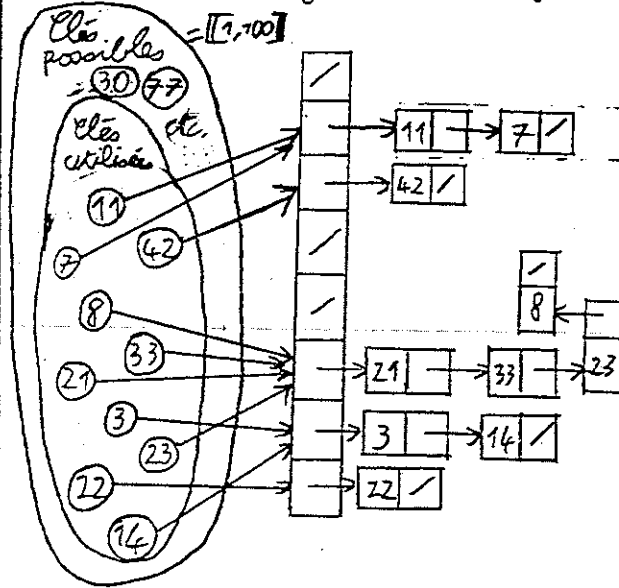
	Recherche	Insertion	Suppression	Recherche par rang
Liste (non triée)	$O(n)$	$O(1)$	$O(n)$	$O(n)$ (ou $O(n^2)$)
Tableau trié	$O(\log n)$	$O(n)$	$O(n)$	$O(1)$
ABR (pire cas)	$O(n)$	$O(n)$	$O(n)$	$O(n)$
ABR (moyenne) = AVL	$O(\log n)$	$O(\log n)$	$O(\log n)$	$O(n)$
Hachage (pire)	$O(n)$	$O(n)$	$O(n)$	$O(n)$
Hachage (moy)	$O(1)$	$O(1)$	$O(1)$	$O(n)$

[Frisolavaux]

Annexe 1: Arbre binaire



Annexe 3: Hashage avec chaînage

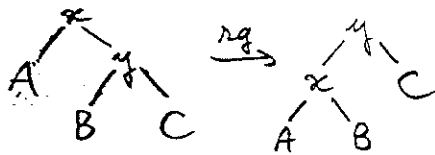


Références:

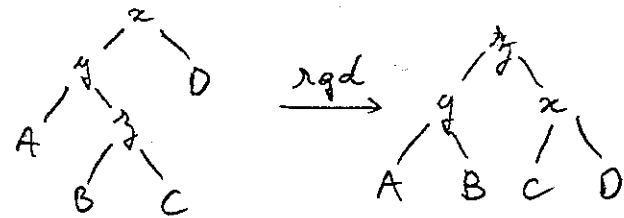
- [Cor] Cormen et al.
- [PS] Preparata-Shamos, Computational Geometry - an Introduction. (Beaulieu)
- [Sed] Sedgewick, Algorithms (Beaulieu)
- [Froi] Froidevaux-Gaudel-Soria Types de données et algorithmes
- [Knu] Knuth, The art of computer programming, Part 2: Sorting and Searching

Annexe 2: Rotation d'arbres

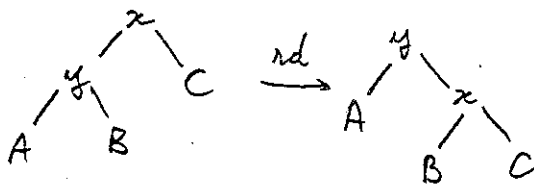
Gauche:



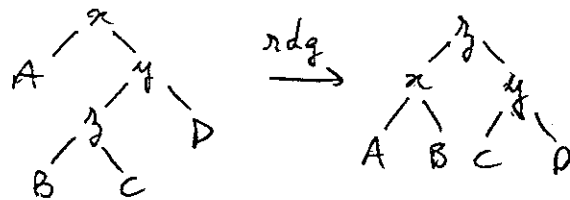
Gauche-droite:



Droite:



Droite-gauche:



Il y a pas eu d'apps
 ds le titre.
 - P. è en dipl, bof

Intersection de segments dans le plan.

[Preparata - Shamos]
 Computational geometry,
 an introduction
 7.2 p. 281 et suivantes

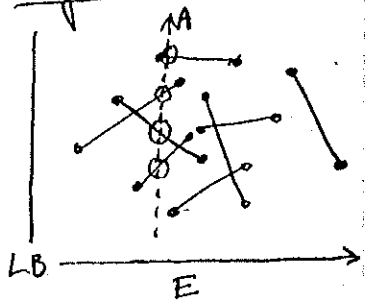
(Algo de Bentley-Ottmann) ↓

Le problème: + TROP long - Faire avec des segments horizontaux et verticaux?
Entrée: N segments du plan, donnés par les coordonnées de leurs extrémités, que l'on suppose toutes distinctes, ainsi que leurs intersections.
Sortie: La liste de leurs points d'intersection. celles de

Algorithme naïf: On cherche les points d'intersection pour chaque paire de segments: $O(N^2)$; c'est même la complexité optimale, puisque sur une grille $\left\{ \begin{array}{c} \text{---} \\ \text{---} \\ \text{---} \\ \text{---} \\ \text{---} \end{array} \right\}_{N/2}$, il y a $O(N^2)$ points d'intersection.

On présente un algorithme en $O((N+I) \log_2 N)$ où I est le nombre de points d'intersection, ce qui en fait un algorithme ^{plus} efficace si $I = o(N^2)$.

Algorithme de balayage:



On balaye la figure avec une ligne de balayage verticale LB de gauche à droite.

On maintient un AVL E contenant les événements qui indiquent les points d'importance pour LB: extrémités (données a priori) et intersections (déterminées pendant le balayage).

Pour détecter les intersections, on constate que deux segments qui se croisent doivent être "adjacents" immédiatement avant, c'est-à-dire que leurs intersections avec LB doivent être côte-à-côte si on les trie par ordonnées croissantes.

Or, l'adjacence de deux segments (relativement à LB) n'est modifiée que lorsqu'un segment commence ou termine, ou lorsqu'un point d'intersection est atteint.

On maintient donc un AVL A pour avoir accès aux segments qui intersectent LB à un moment donné, dans l'ordre d'ordonnées croissantes.

* Initialisation: On stocke les coordonnées des N segments dans un AVL E en utilisant l'abscisse comme clé. On initialise l'AVL A à \emptyset , et on pose $L =$ liste vide (liste des intersections détectées) dans un tour de boucle.

* Boucle: Tant que E n'est pas vide, soit $e = \min E, L = \emptyset$.

① Si e est un début de segment s .

Inserer s A

$s_1 = \text{succ } s$ A

$s_2 = \text{préd } s$ A

) nouvelles adjacences

Si $s_1 \cap s_2 \neq \emptyset$, alors $L \leftarrow (s_1, s_2)$ Notation nd

② Si e fin de segment s

$s_1 = \text{succ } s$ A

$s_2 = \text{préd } s$ A

) nouvelles adj

Si $s_1 \cap s_2 \neq \emptyset$ à droite de e , alors $L \leftarrow (s_1, s_2)$

Supprimer s A

③ Si $e = s_1 \cap s_2$ (avec $s_1 > s_2$ dans A avant e)

$s_3 = \text{succ } s_1$ A

$s_4 = \text{préd } s_2$ A

) nouv. adj

Si $s_3 \cap s_4 \neq \emptyset$ alors $L \leftarrow (s_3, s_4)$
 (s_4, s_1)

Echanger s_1, s_2 dans A .

Afficher e

④ Pour $(s_1, s_2) \in L$

Si $s_1 \cap s_2 \notin E$, alors Inserer $s_1 \cap s_2$ E

↑
Rechercher $s_1 \cap s_2$ E

Il manque l'affichage.

Correction: Chaque intersection est détectée une fois (parce que toutes les adjacences sont mises à jour) et une seule (grâce à ④).

Complexité: Initialisation: $O(N \log N)$

①, ② et ③: $O(\log N)$ (opérations dans A)

④: On a détecté au plus 2 intersections, donc $O(\log(N+K)) = O(\log N)$

Boucle: $O((N+K) \log N)$

⇒ $O((N+K) \log N)$

Simplifier pour détailler les structures de données.

Hasage Parfait

Thm: On peut stocker n clés (statiques) dans une structure de données donnant une recherche en $O(1)$ et de taille mémoire $O(n)$

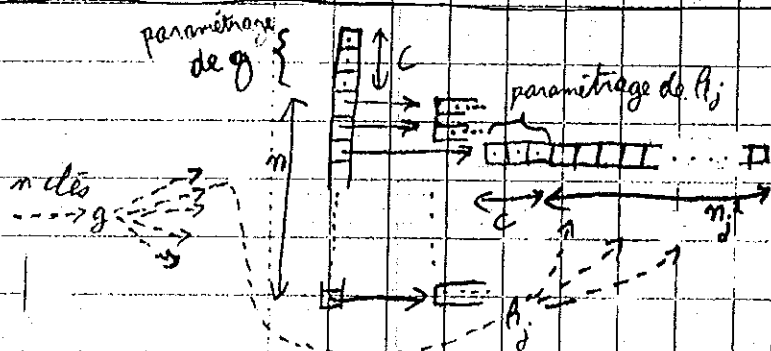
préreqs:
hasage uniforme
classe universelle.

Idée:

- Construire un hasage sur une mémoire $O(n^2)$ sans collisions (\rightarrow recherche en $O(1)$)
- Faire un "grand hasage" à n alvéoles, puis utiliser le hasage précédent sur chacune des alvéoles

Lemme 1: Si on stocke n clés dans une table de taille n^2 avec une fonction h de hasage uniforme choisie aléatoirement, alors la probabilité d'avoir une collision est inférieure à $\frac{1}{2}$

Lemme 2: Si on stocke n clés dans une table de taille n^2 avec une fonction g de hasage uniforme (aléatoire), alors, en notant n_j le nombre de clés envoyées sur l'alvéole j ($1 \leq j \leq n$), on a $E[\sum_{j=1}^n n_j^2] < 2n$ et donc $P(\sum_{j=1}^n n_j^2 \geq 4n) < \frac{1}{2}$



Espace total nécessaire:
 $C + n + nC + \sum_{i=1}^n n_i = O(n)$
 (quitte à relâcher des choix pour g et les h_i)

(Le C vient du fait que l'on peut paramétrer un choix dans une classe universelle par trois ou quatre constantes.)

Démo Lemme 1:

Un couple de lés étant donné, il y a une probabilité de $\frac{1}{n^2}$ qu'ils soient envoyés sur la même valeur.

sachant qu'il y a $\binom{n}{2}$ lés, on a; \rightarrow avec un tirage uniforme
avec X la variable aléatoire comptant le nombre de collisions,

$$E(X) = \binom{n}{2} \cdot \frac{1}{n^2} = \frac{1}{2} - \frac{1}{2n} < \frac{1}{2}$$

D'où $P(X \geq 1) \stackrel{(*)}{\leq} E(X) < \frac{1}{2}$. on a moins d'une chance sur deux d'avoir une collision. \square

Démo Lemme 2:

On remarque: $\forall a \in \mathbb{N}, a^2 = a + 2 \binom{a}{2}$

$$\begin{aligned} \text{Alors } E\left[\sum_{j=1}^n n_j^2\right] &= E\left[\sum_{j=1}^n n_j\right] + 2E\left[\sum_{j=1}^n \binom{n_j}{2}\right] \rightarrow \text{linéarité de l'espérance} \\ &= n + 2E\left[\sum_{j=1}^n \binom{n_j}{2}\right] \rightarrow \sum n_j = n \text{ et } n_j \text{ sont pas aléatoires} \end{aligned}$$

de plus, $\sum_{j=1}^n \binom{n_j}{2}$ est le nombre total de collisions

Comme g est un tirage uniforme, cette somme est bornée par $\binom{n}{2} \frac{1}{n} = \frac{n-1}{2}$

$$\text{donc } E\left[\sum_{j=1}^n n_j^2\right] \leq n + 2 \frac{n-1}{2} < 2n$$

$$\text{On a alors } P\left(\sum_{j=1}^n n_j^2 \geq 4n\right) \stackrel{(*)}{\leq} \frac{E\left[\sum_{j=1}^n n_j^2\right]}{4n} < \frac{1}{2} \square$$

(*) on utilise l'inégalité de Markov: $P(X \geq c) \leq \frac{E(X)}{c}$