

À partir d'un ensemble de données de même type, on va chercher des structures de données adaptées pour effectuer les opérations :

- insérer.
- supprimer.
- rechercher.

I Méthodes simples.

A) Recherche dans une liste non triée.

Pour rechercher un élément dans la liste : on le compare avec le premier élément ; si c'est négatif, on passe au suivant... Sur exemple, la suppression infructueuse est en $O(n)$.

B) Recherche dans une liste triée.

On a une relation d'ordre sur les éléments. L'insertion d'un élément en fin de liste se fait en $O(n)$ comparaisons. Performance moyenne : $O(n)$.

C) Recherche dichotomique.

Recherche d'une occurrence d'un élément quelconque se fait en $O(\log n)$ comparaisons. On utilise du récursif.

En bilan : un algorithme fonctionnant par comparaisons ne peut faire mieux que $O(\log n)$ (arbre de décision).

II Structure d'arbre.

A) Arbres de recherche.

Def 1: Un univers U est un ensemble totalement ordonné. Pour éviter de déplacer les données $x = \{x, y, z\}$ à l'hier, on associe à chaque donnée une clé, qui est un élément de l'univers. Réciproquement, une clé nous permet d'accéder à une unique donnée : de $c(x)$ on trouve x .

Def 2: Un arbre binaire de recherche est un arbre muni d'une fonction clé : $\forall x$ sommet, $\forall y \in A_g(x)$ (sous-arbre gauche du nœud x)
 $\forall z \in A_d(x)$ (sous-arbre droit du nœud x),
 $c(y) < c(x) < c(z)$.

Def 3: Une arborescence ordonnée de recherche est un arbre tel qu'en chaque nœud x , à $d(x)$ fils et muni de $d(x)-1$ clés $c_1(x), \dots, c_{d(x)-1}(x)$, envoyant $A_1(x), \dots, A_{d(x)}(x)$ les sous-arbres de x , pour tout y_i sommet de $A_i(x)$, on a :
 $c(y_1) < c_1(x) < c(y_2) < \dots < c_{d(x)-1}(x) < c(y_{d(x)})$.

Def 4: Un arbre de recherche est un ABR ou un AOR.

Rem 5: Ces structures sont dynamiques, elles sont amenées à changer au cours du temps (exemple: dictionnaire).

Def 6: On effectue les opérations suivantes sur les arbres binaires : rotations gauche, droite, (quad) q-d ou d-q. (voir annexes 1-2).

Prop 7: Si A est un ABR, alors les rotations sur A préservent cette structure. De plus, elles se réalisent en temps constant.

B) Arbres AVL

Def 8: On note S la fonction "différence de hauteurs" définie sur l'ensemble des ABR : $S(A) = \text{hauteur}(A_g) - \text{hauteur}(A_d)$.

Les AVL sont les ABR respectant la condition : $S(A) \in \{-1, 0, 1\}$.

Prop 9: Soit A un AVL ayant n sommets et de hauteur h . Alors :

$$\log_2(1+n) \leq 1+h \leq 1,44 \log_2(2+n)$$

Cette borne minimale est essentiellement atteinte par les arbres de Fibonacci :

$$F_0 = 0, F_1 = 1, F_{n+2} = F_{n+1} + F_n$$

Prop 10: Les opérations insertion, suppression et recherche se font en temps $O(\ln n)$. [DEVI]

C) Arbres a-b.

Def 11: Soient a et b deux entiers avec $a \geq 2$ et $b \geq 2a-1$. Un arbre $a-b$ est un arbre A vérifiant :

- i) toutes les feuilles ont la même profondeur.
- ii) la racine a au moins 2 et au plus b fils.
- iii) les autres nœuds ont au moins a et au plus b fils.

Soit $d(x)$ le nombre de fils de x , et $A_i(x)$ le i -ème sous-arbre de x .

Def 12: Soit S un ensemble de clés. A est un arbre $a-b$ pour S si les éléments de S sont rangés aux feuilles de A en ordre croissant de la gauche vers la droite, et si pour chaque nœud x à $d(x)-1$ clés $k_1 < \dots < k_{d(x)-1}$ (les balises de x) vérifient :

pour toute clé c_i de $A_i(x)$, on a :

$$c_1 \leq k_1 < c_2 \leq k_2 < \dots \leq k_{d(x)-1} < c_{d(x)}$$

26/07
2015

Q 21 : Algorithme de recherche et structures de données associées.

[FR0]

[BBC p 147]

[BBC p 151]

[BBC p 160]

Ex 13: voir annexe 3.

Prop 14: La recherche dans un arbre a-b se fait en temps $O(\log n)$.

Principe 15 de l'insertion: ① On descend dans l'arbre à partir de la racine, et on ajoute au père une balise (cf annexe 4).

② On vérifie que le père a entre a et b fils. Si ce n'est pas le cas, on réalise un éclatement (cf annexe 5), et si nécessaire, on continue jusqu'à la racine.

Prop 16: L'insertion dans un arbre a-b se fait en temps $O(\log n)$.

Principe 17 de la suppression: on recherche la clé, on la supprime, on supprime une balise (si nécessaire, réaliser une fusion).

Prop 18: La suppression dans un arbre a-b se fait en temps $O(\log n)$.

D) Arbres bicolorés (rouge-noir)

Def 19: Un arbre bicoloré (a, c) où A est un arbre complet binaire et c une coloration des sommets dans {rouge, noir} est un arbre vérifiant:

- a) toutes les feuilles sont noires;
- b) la racine est noire;
- c) le père d'un sommet rouge est noir;
- d) les chemins d'un même sommet à une feuille ont le même nombre de sommets noirs.

Un arbre bicoloré pour un ensemble ordonné S est un arbre bicoloré qui est un ABR pour S, où les éléments de S ne sont que ses nœuds.

Ex 20: voir annexe 6:

Def 21: Une fonction rang rg définie sur un ABR complet A est une application de S (sommets de A) dans \mathbb{N} vérifiant:

- i) \forall feuille x , $rg(x) = 0$ et si x a un père, $rg(p(x)) = 1$.
- ii) \forall sommet x , $rg(x) \leq rg(p(x)) \leq rg(x) + 1$.
- iii) \forall sommet x , $rg(x) < rg(p(p(x)))$.

Prop 22: si (A, c) est un arbre bicoloré, alors il existe une fonction rang sur A. Réciproquement, l'existence d'une fonction rang indique qu'il existe

une coloration pour A tel qu'il soit bicoloré.

Prop 23: les opérations d'insertion et de suppression se font en temps $O(\log n)$.

E) Arbres binaires de recherche optimisés.

Def 24: Soient h_1, \dots, h_n n des trices par ordre croissant. La probabilité qu'une clé h_i soit demandée est p_i . On ajoute $n+1$ des factices, d_0, \dots, d_n , où d_i représente des mots recherchés entre h_i et h_{i+1} . Soit q_i la probabilité qu'on recherche d_i . Un ABRO est un arbre binaire minimisant:

$$|E| \text{ [coût de recherche dans A]} = 1 + \sum_{i=1}^n \text{prof}_A(h_i) p_i + \sum_{i=0}^n \text{prof}_A(d_i) q_i$$

Rem 25: Cette structure est statique, l'insertion et la suppression d'un élément entraînent une reconstruction totale de l'arbre.

Ex 26: cf annexe 7.

IV Application à l'algorithmique du texte.

A) Plus longue sous-séquence commune.

Exemple: Comparaison de deux liens d'ADN.

Def 27: $Z = \langle z_1, \dots, z_h \rangle$ est une sous-séquence de $X = \langle x_1, \dots, x_m \rangle$,

si il existe $\varphi: [1; h] \rightarrow [1; m]$ strictement croissante telle que

$$\forall i \in [1; h], z_i = x_{\varphi(i)}$$

On recherche la plus longue sous-séquence commune (PLSC) à

$X = \langle x_1, \dots, x_m \rangle$ et à $Y = \langle y_1, \dots, y_n \rangle$.

Prop 28: Soit Z une plus longue sous-séquence commune à X et à Y.

Si $z_h = x_m = y_n$, alors Z_{h-1} est une PLSC à X_{m-1} et Y_{n-1} .

Si $z_h \neq x_m$ alors Z_h est une PLSC à X_{m-1} et Y_n .

Si $z_h \neq y_n$ alors Z_h est une PLSC à X_m et Y_{n-1} .

Résolution par programmation dynamique; récurrence sur la taille de X et sur la taille de Y.

B) Distance d'édition.

Ex: Comparer les mots REMER et RUITEUR.

[ABC]

[COR p 365]

[COR p 362]

On veut minimiser la distance d'édition avec trois possibilités d'action,

- aligner une lettre de x et une de y .
- sauter une lettre de x .
- sauter une lettre de y .

ex: $REMUE - ER$
 $R - - VMEUR$
 $\uparrow \quad \uparrow \quad \uparrow$
 distance est de 4.

$O(nm)$

c) Tri des suffixes.

Pour rechercher un mot dans un texte : création de la table des suffixes. Un mot est défini par une suite de lettres $y[0], \dots, y[n-1]$, et on veut lister les suffixes $y[0, \dots, n-1]$, $y[1, \dots, n-1]$, ..., $y[n-1]$.

On va chercher une permutation p telle que $y[p(0), \dots, n-1] < y[p(1), \dots, n-1] < \dots$.

Def 29: Soit $h > 0$. Soit $u \in A^*$ le début d'ordre h du mot u :

$$\text{prem}_h(u) = \begin{cases} u & \text{si } |u| \leq h \\ u[0, \dots, h-1] & \text{sinon.} \end{cases}$$

On note $R_h[i]$ le rang de $\text{prem}_h(y[i, \dots, n-1])$ dans l'ordre croissant des $(R_h[i])_{i=0, \dots, n-1}$. Pour tout h , on a la relation d'équivalence: $i \equiv_h j$ si $R_h[i] = R_h[j]$.

Soit $i \geq n$, on pose $R_h[i] = -1$.

emme 30 (de dédoublement): $R_{2h}[i]$ est le rang du couple

$(R_h[i], R_h[i+h])$ dans la liste croissante de ces couples.

Th 31: L'algorithme TRI-SUFFIXES appliqué à y de longueur n crée la table des suffixes en temps $O(n \log n)$

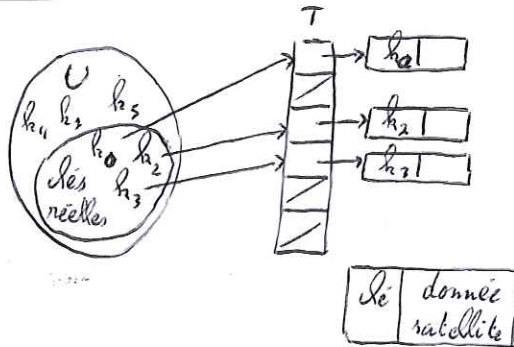
[DVP]

IV Tables de hachage.

A) Tables à adressage direct.

Def 32: Une table à adressage direct est un tableau $T[0, \dots, m-1]$ dans lequel chaque position (ou adresse) correspond à une clé de l'univers U .

Pr op 33: Chacune des opérations : recherche, insertion et suppression se fait en temps $O(1)$.



B) Tables de hachage.

Problème: Si $|U|$ grand, on ne peut pas toujours gérer T aussi grand.

Def 34: Une fonction de hachage est une fonction $h: U \rightarrow [0, m-1]$, où $|m| \ll |U|$; le tableau $T[0, \dots, m-1]$ est la table de hachage. On dit qu'un élément de U est haché dans l'adresse $h(k)$.

Inconvénient: la non injectivité de h peut envoyer deux clés dans la même adresse: on a un phénomène de collision! Pour le résoudre, on utilise le chaînage (cf annexe 8).

Def 35: On définit $\alpha = \frac{n}{m}$ le facteur de remplissage, c'est-à-dire pour n éléments stockés dans un tableau de taille m .

Th 36: Dans une table de hachage dont la résolution se fait par chaînage, le coût d'une recherche réussie ou d'une recherche infructueuse est $O(1 + \alpha)$, dans le cas d'un hachage uniforme simple (chaque élément a la même chance d'être haché vers l'une quelconque des adresses).

c) Des exemples de fonctions de hachage.

Def 37: Méthode de la division (les clés sont des entiers).

$$h: k \mapsto h \text{ [mod } m \text{]}.$$

Si les données sont stockées en bit, on évite que m soit $m = 2^p$.

Def 38: Méthode de la multiplication. Soit $0 < A < 1$.

$$h: k \mapsto m(hA - [hA]) \text{ [mod } 1 \text{]}.$$

Def 39: Si la fonction de hachage est donnée, un ennemi peut choisir les clés telles qu'elles soient stockées dans la même cellule. On sélectionne donc aléatoirement une fonction de hachage: c'est le hachage universel.

Def 40: Une collection de fonctions de hachage est universelle si:

$$\# \{h \in \mathcal{H}; h(k) = h(l)\} \leq \frac{|\mathcal{H}|}{m}, \quad \forall k, l \in U.$$

Th: Soit n_i la longueur de la liste $T[i]$. Soit une fonction de hachage choisie aléatoirement dans une collection universelle, et k la clé recherchée: si k est dans la table: $\mathbb{E}[n_i(h(k))] \leq \alpha = \frac{n}{m}$; sinon: $\mathbb{E}[n_i(h(k))] \leq \alpha + 1$.

