

Algorithmes de recherche et structure de données associées.

921

Introduction :

* Motivation : Stocker et exploiter des données est une utilisation des plus communes de l'informatique. Savoir retrouver une donnée précise dans une collection est un problème récurrent.

* Problème de recherche :

entrée : K une collection d'objets ; $E \subseteq K$ et $R \in K$.

sortie : $R \in K$.

I. Recherche dans un dictionnaire.

Objectif : Recherche d'un élément quand $K = \{(cl, valeur)\}$ et E est un dictionnaire (structure abstraite).

A. L'opération recherche suit une loi uniforme sur les clés.

Objectif : Recherche d'un élément quand $K = \{(cl, valeur)\}$ et E est un dictionnaire. De plus, R est uniformément distribué sur E .

Algorithme 1 : Recherche naïve : on compare R à tous les éléments de E .

- Structures de données associées : liste, tableau
- Complexité temporelle dans le pire cas : $O(m)$.

Algorithme 2 : Recherche dichotomique. [1, p. 175]

- Structure de données associée : tableau trié

dicho(R, E, g, d, res):

$g, d \in \{0, \dots, E -1\}$
$res \in \{0, \dots, E \}$

```

Si  $g \leq d$  alors
   $m := (g+d) \div 2$ 
  Si  $R = E[m]$  alors  $res := m$ 
  Sinon Si  $R < E[m]$ 
    alors dicho( $R, E, g, m-1, res$ )
  Sinon dicho( $R, E, m+1, d, res$ )
Sinon  $res := 0$ 
  
```

→ Complexité temporelle dans le pire cas : $O(\log m)$

Algorithme 3 : Recherche supportée par un ABR équilibré : recherche dans un ABR.

Définition 4 : [1, p. 197] Un arbre binaire de recherche (ABR) est un arbre binaire étiqueté tel que pour tout nœud x d'étiquette m de l'arbre :

- * les éléments de tous les nœuds du fils gauche de x sont inférieurs à m ;
- * les éléments de tous les nœuds du fils droit de x sont supérieurs à m .

(Figures en annexe : 1 et 2)

Définition 5 : [1, p. 221] Un ABR à m nœuds est dit équilibré si sa hauteur est en $O(\log m)$.

→ Implémentations : AVL [1, p. 224]

arbre rouge-noir [2, p. 287]

→ Complexité temporelle dans le pire cas : $O(\log m)$

Application 6 : Recherche dans le cache de la géométrie algorithmique.

→ Problème : L'intersection de m droites par la méthode de balayage.

On cherche un élément dans K est l'ensemble des points du plan et E l'ensemble des points d'intersection de deux droites.

• Implémentation : arbre rouge-noir

• Complexité temporelle dans le pire cas : $O(m \log m)$.

→ Problème : Recherche d'un point dans un nuage.

On cherche un élément dans $K = \{(cl, valeur)\}$ avec cl dans l'ensemble des points du plan et E est l'ensemble des points du nuage.

• Solution : on utilise les quad-trees

(Figure en annexe : 4)

[2, p. 940]

[3, p. 309]

B - La recherche n'est pas uniforme sur l'ensemble des clés.

Objectif: Recherche d'un élément quand $K = \{(clé, valeur)\}$ et E est un dictionnaire et P suit une distribution sur E .

Algorithme 7: Recherche auto-adaptative [1, p. 175]

- On met en première position l'élément que l'on vient de rechercher.
- Structure de données associée: liste
- Complexité temporelle asymptotique en le nombre de recherche et en moyenne: $O(\frac{n}{P_{eq} m})$ [4]

Algorithme 8: Recherche supportée par un ABR optimal: recherche dans un ABR [2, p. 368]

Définition 9: Un ABR optimal minimise la complexité espérée de la recherche en fonction de la distribution sur les clés.

(Figure en annexe: 5)

→ Implémentation: arbre splay.

C - Les clés ne possèdent pas nécessairement un ordre total. [2, p. 239]

Objectif: Recherche d'un élément quand $K = \{(clé, valeur)\}$ où les clés sont à valeur dans un ensemble qui n'admet pas nécessairement un ordre total.

Principe 10: Pour chaque valeur de clé, une fonction de hachage h lui associe un entier représentant un indice de tableau qui correspond à la case qui va la stocker.

Définition 11: Deux éléments $u \neq v$ sont en collision si $h(u) = h(v)$.

- Implémentation: table de hachage.
 - Gestion des collisions par chaînage simple.
 - Complexité temporelle dans le pire cas: $O(n)$

Définition 12: On parle de hachage parfait si la recherche s'effectue en $O(1)$ (Figure en annexe 6)

D - Le dictionnaire doit être stocké sur un disque externe
Objectif: Recherche d'un élément quand $K = \{(clé, valeur)\}$ et E est un dictionnaire qui ne peut pas être stocké en mémoire vive.

Définition 13: Un arbre B est un arbre de recherche qui minimise le nombre d'accès au disque externe. [DEV]

Remarque 14: Un arbre B peut être étendu en un arbre B^+ pour faciliter les recherches séquentielles.

Application 15: Base de données relationnelle. (Figure en annexe 7; 8)

III - Recherche dans un texte.

Objectif: Recherche d'un élément lorsque K est l'ensemble des mots de Σ^* (Σ alphabet) et E l'ensemble des facteurs d'un texte.

Application 16: Recherche d'un mot dans un fichier texte, d'une chaîne de caractère dans une séquence ADN, ...
 → Structure de données associée: tableau de la taille de la chaîne.

Algorithme 17: Recherche naïve. [2, 908]

Recherche naïve (T, P)
 $m := T$. longueur; $n := P$. longueur
 pour $s = 0$ à $n - m$
 si $P[1, \dots, m] = T[s+1, \dots, s+m]$
 alors retourner s

entree: T, texte
 P, motif
 sortie: do'ralage

→ Complexité temporelle dans le pire cas: $O((n - m + 1)m)$

A - Algorithme de Rabin-Karp [2, 910]

Hypothèse: On considère que $\Sigma = \{0, \dots, 9\}$.

Algorithme 18: Algorithme de Rabin-Karp

Rabin-Karp (T, P)
 $m := T$. longueur; $n := P$. longueur; $R := d^{m-1} \text{ mod } q$
 $p := 0$; $t_0 := 0$
 pour $i = 1$ à m
 $p := (dp + P[i]) \text{ mod } q$; $t_0 := (dt_0 + T[i]) \text{ mod } q$
 pour $s = 0$ à $n - m$
 si $p = t_0$ alors
 si $P[s+1, \dots, s+m] = T[s+1, \dots, s+m]$ alors retourner s

entree: T, texte
 P, motif
 sortie: do'ralage

$t_{s+1} = (d(t_0 - T[s+1]R) + T[s+m+1]) \text{ mod } q$

[2, p. 239]

[2, p. 239]

→ Complexité temporelle du prétraitement : $O(m)$
 → Complexité temporelle dans le pire cas : $O((n-m+1)m)$.
 (Figures en annexe 9)

B. Algorithme de Knuth-Morris-Pratt [2, p. 215]

Hypothèse : Le motif est fixe et connue à l'avance.

Définition 19 : La fonction suffixe associée au motif P , $\sigma : \Sigma^* \rightarrow \{0, \dots, m\}$ donne la longueur du plus long préfixe de P qui est suffixe de $x \in \Sigma^*$.

Définition 20 : On définit l'automate des occurrences associé à $P \in \Sigma^+$ comme suit : $Q = \{1, \dots, m\}$; $q_0 = 0$; $q_f = m$; $\delta(q, a) = \sigma(P_q a)$ pour $a \in \Sigma$, $q \in Q$ et $P_q = P[1..q]$

Algorithme 21 : Calcul de S .

calcul - $S(P, \Sigma)$

```

  m := P.longueur
  pour q = 0 à m
    pour a ∈ Σ
      R := min(m+1, q+2)
      répéter
        R := R-1
      jusqu'à P_R suffixe de P_q a
      S(q, a) := R
  retourner S
  
```

entrée : P motif
 Σ alphabet
 sortie : S

→ Complexité temporelle du prétraitement : $O(m|\Sigma|)$

Algorithme 22 : Recherche dans l'automate

Recherche - automate (T, δ, m)

```

  m := T.longueur ; q := 0
  pour i = 1 à m
    q := δ(q, T[i])
    si q = m
      retourner i-m.
  
```

entrée : T texte
 δ : fonction de transition
 m : longueur du motif
 sortie : décalage

→ Complexité temporelle dans le pire cas : $O(m)$

Définition 23 : La fonction préfixe associée au motif P , $\pi : \{1, \dots, m\} \rightarrow \{0, \dots, m-1\}$ donne la longueur du plus long préfixe de P qui est suffixe de P_q où $q \in \{1, \dots, m\}$.

Algorithme 24 : Algorithme de Knuth-Morris-Pratt

→ Idée : construction intelligente de l'automate précédent pour que la fonction S puisse être calculée à la volée.

→ Complexité temporelle du prétraitement dans le pire cas : $O(m)$
 → Complexité temporelle dans le pire cas : $O(m)$ **DEV**

IV. Exercices

Problème 25 : Problème de la recherche du min et du max :

On cherche un élément k , dans K un ensemble ordonné et $E \subseteq K$, $R = \inf E$ ou $k = \sup E$

Problème 26 : Union - find et partition d'un ensemble :

On cherche un élément k , dans K l'ensemble des entiers et $E = \{(m, \bar{m})\}$ où \bar{m} représente son représentant de sa classe d'équivalence.

Problème 27 : Problème de la recherche non-associative :

on ne recherche pas sur la valeur de k mais également sur son contexte : recherche "temporelle" (l'élément qui est dans la structure depuis le plus longtemps), de K élément.

Problème 28 : Problème d'accessibilité dans un graphe.

On cherche un élément k , dans K l'ensemble des états d'un graphe et E l'ensemble des états accessibles de ce graphe.

References:

- [1] Fraedensauz, Gaudel, Types de données et algorithmes
- [2] Rivest, Stein, Cormen, Leiserson, Algorithmique, 3^{ème} édition.
- [3] Overmars, Schwarzkopf, Berg, van Kreveld, Computational Geometry
- [4] Knuth, The art of Computer Programming, vol 3

Exercice:

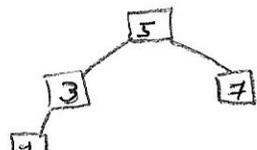


Figure 1: Un ABR équilibré

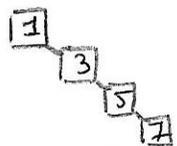


Figure 2: Un ABR déséquilibré

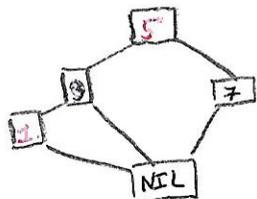


Figure 3: Un arbre rouge-noir

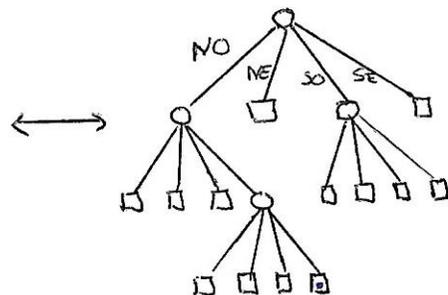
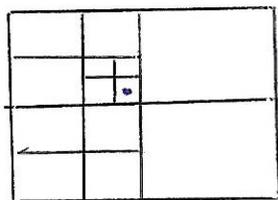
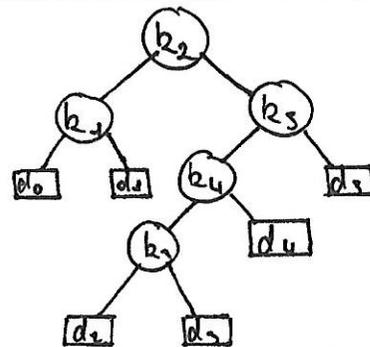


Figure 4: Un exemple de Quad-Tree.



i	0	1	2	3	4	5
P _i		0,25	0,1	0,05	0,1	0,2
q _i	0,05	0,1	0,05	0,05	0,05	0,1

Figure 5: Un exemple d'ABR optimal.

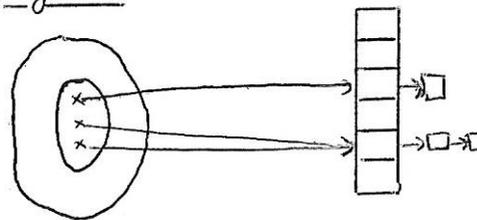


Figure 6: Hachage avec résolution des collisions par chaînage simple.

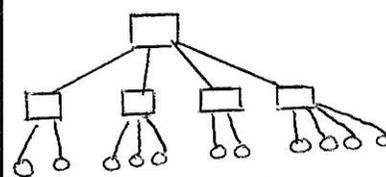


Figure 7: Un arbre B

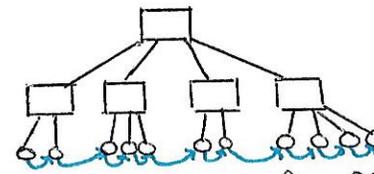
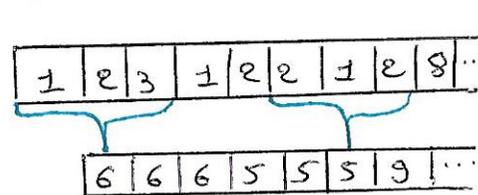


Figure 8: Un arbre B+



On cherche P = [2 | 1 | 2]

Figure 9: Algorithme de Rabin-Karp.