

EA4 – Éléments d’algorithmique

TD n° 11 : hachage (suite)

Exercice 1 : tables de hachage et adressage ouvert

On considère une table de hachage T à adressage ouvert, de longueur $m = 11$ avec la fonction de hachage $h(k) = k \bmod 11$.

- Insérer successivement dans T les clés 5, 28, 19, 15, 20, 33, 37, 17, 26 et 10 en appliquant les méthodes de résolution des collisions suivantes :
 - sondage linéaire
 - double hachage avec $h_1(k) = h(k)$ et $h_2(k) = 1 + (k \bmod (m - 1))$.
- Supprimer ensuite les clés 19, 37 et 17 en détaillant la suite de sondages réalisés.
- Comment se passe ensuite l’insertion de 26 ? de 39 ? Préciser quels sont les sondages réalisés et la case finalement choisie pour l’insertion.
- Plus généralement, quel est le bon paramètre à considérer pour décider quand procéder à un « redimensionnement » (qui peut éventuellement se borner à un simple nettoyage sans changement de dimension) ? Est-ce le même pour les insertions et les suppressions ?

Exercice 2 : hachage et doublons

Soit T un tableau de taille n possédant au plus m entiers différents. On veut calculer son nombre exact de valeurs différentes.

- Donner une solution en $O(n \log n)$ ne faisant pas appel au hachage.
- Donner une solution faisant appel au hachage. Quelle sera sa complexité en mémoire ?
- Avec l’hypothèse simplificatrice que l’on a accès à une fonction de hachage sans collision sur nos données (c’est-à-dire telle que $h(x) \neq h(y)$ pour deux éléments distincts dès lors que la table est assez grande), quelle est la complexité de l’algorithme ?
- Quelle est la complexité dans le pire cas si l’on suppose qu’on n’a jamais plus que des k -collisions (c’est-à-dire que jamais plus de k valeurs différentes ont la même image par la fonction de hachage) ?
- Sous quelles hypothèses moins restrictives a-t-on une complexité *moyenne* du même ordre ?

Exercice 3 : borne en moyenne de plus long sondage

Le but de cet exercice est de prouver que dans une table à adressage ouvert, avec taux de remplissage α borné et hachage supposé uniforme, la longueur moyenne de la plus longue suite de sondages est en $O(\log n)$.

On se place pour simplifier dans le cas $\alpha = \frac{1}{2}$, avec table initialement vide de longueur au moins $2n$, et on insère successivement n éléments. On note S_i le nombre de sondages nécessaires pour la i^{e} insertion.

- Rappeler pourquoi $\mathbb{P}(S_i > k) \leq 1/2^k$ pour tous i et k .
- En déduire que $\mathbb{P}(S_i > 2 \log n) \leq 1/n^2$ pour tout i .
- Soit $S = \max_{1 \leq i \leq n} S_i$. Montrer que $\mathbb{P}(S > 2 \log n) \leq 1/n$.
- En déduire que $\mathbb{E}(S) \in O(\log n)$.

Exercice 4 : algorithme de Rabin et Karp (examen 2015)

L'objectif de cet exercice est de donner une application des fonctions de hachage à la recherche d'un motif dans un texte.

On commence par proposer une méthode naïve.

1. Écrire une fonction `test_occurrence(m, t, i)` retournant `True` si un motif `m` apparaît à la position `i` d'un texte `t` (`m` et `t` étant représentés par des tableaux de caractères).
2. Écrire une fonction `recherche_naive(m, t)` effectuant une recherche naïve du motif `m` dans `t` à l'aide de `test_occurrence(m, t, i)`.
3. Quelle(s) opération(s) est-il raisonnable de prendre en compte pour le calcul de la complexité (en temps) des algorithmes précédents? En déduire la complexité de l'algorithme `recherche_naive(m, t)`. Justifier.

Dans la suite on fera l'hypothèse que les chaînes de caractères ont été remplacées par des tableaux d'entiers en considérant la valeur entière du code binaire de chaque caractère (via la fonction `ord()` en Python, ou une conversion de format en Java).

L'algorithme de Rabin et Karp améliore l'algorithme `recherche_naive(m, t)` précédent. Il repose sur l'utilisation d'une fonction de hachage.

Soit $b \geq 2$ un entier fixé ; pour tout entier $k \geq 1$, on considère la fonction $h_k : \mathbb{N}^k \rightarrow \mathbb{N}$ telle que :

$$h_k(x_{k-1}, \dots, x_0) = x_{k-1} \cdot b^{k-1} + x_{k-2} \cdot b^{k-2} + \dots + x_2 \cdot b^2 + x_1 \cdot b + x_0$$

4. Rappeler le principe de la méthode de Horner pour évaluer un polynôme.
5. Soit `t` un tableau d'entiers, de longueur $\ell \geq k$. Décrire un algorithme `initialise_h(t, k)` permettant de calculer $h_k(t[0], \dots, t[k-1])$ en $\Theta(k)$ opérations arithmétiques sur des entiers (additions ou multiplications).
6. Si les entiers considérés sont quelconques, ce calcul est-il réellement de complexité linéaire? Et si tous les calculs sont faits *modulo* un entier q fixé?

Dorénavant, tous les entiers sont considérés modulo q (inutile de réécrire `initialise_h(t, k)`).

7. Quelle relation existe-t-il entre $h_k(t[0], \dots, t[k-1])$ et $h_k(t[1], \dots, t[k])$?
8. En déduire un algorithme `affiche_tous_h(t, k)` qui calcule puis affiche *toutes* les valeurs $h_k(t[i], \dots, t[i+k-1])$ pour $i \in [0, \text{len}(t) - k]$ avec une complexité *totale* en $\Theta(\text{len}(t))$.
9. Le principe de l'algorithme de Rabin et Karp est d'utiliser la fonction h_k , pour un k bien choisi, comme un filtre pour limiter le nombre d'appels à `test_occurrence(m, t, i)`. Écrire un algorithme `recherche_RK(t, m)` recherchant le motif `m` dans `t` selon ce principe.
10. Déterminer sa complexité dans le pire des cas, en donnant un exemple de tel pire cas.
11. En faisant l'hypothèse d'uniformité suivante :

b et q peuvent être (et sont) choisis de telle manière que la répartition des images des k -uplets de lettres soit (quasi)-uniforme sur $[0, q-1]$

déterminer la complexité en moyenne de `recherche_RK(t, m)` en fonction de $\ell := \text{len}(t)$, de $k = \text{len}(m)$, de q et du nombre p d'occurrences trouvées. En déduire que, sous des hypothèses raisonnables, cette complexité moyenne est en $\Theta(\ell)$.

Exercice 5 : une application en bioinformatique

En bioinformatique, l'un des domaines principaux de recherche tourne autour de l'analyse de séquences ADN (représentées par des chaînes de caractères sur l'alphabet $\{\text{A}, \text{T}, \text{C}, \text{G}\}$). Par exemple, on peut chercher à repérer un motif codant pour une certaine protéine.

On utilise la fonction de hachage d'un mot $w = w_0 \dots w_n$ définie de la façon suivante :

$$h(w) = \sum_{i=0}^n c_i 4^{n-i} \mod 17, \text{ où } c_i = \begin{cases} 0 & \text{si } w_i = \text{A} \\ 1 & \text{si } w_i = \text{C} \\ 2 & \text{si } w_i = \text{G} \\ 3 & \text{si } w_i = \text{T} \end{cases}$$

1. Calculer $h(\text{CCAGAT})$.
2. Appliquer l'algorithme de Rabin et Karp pour rechercher le mot **CCAGAT** dans le texte
TGCATACGAAATTGCACCAGATCA.