

Leçon 903: Algorithmes de tri: Exemples

07/04
2015

Introduction:

Trier, c'est organiser un ensemble de données suivant un ordre prédéfini.

L'intérêt de leur étude est multiple: améliorer le rendement d'un programme est souvent préalablement. L'entrée, application au domaine graphique (algorithme du peintre)

Parmi les propriétés les plus recherchées, on peut parler de tri -en place, stp modifie directement la structure à trier, sans nécessiter de variables supplémentaires (ou très peu)

-stable, si les éléments égaux pour ce tri sont dans le même ordre à l'entrée et à la sortie

Parmi les tris que nous discuterons, on remarquera d'abord les tris peu comparaisons, avant de finir sur d'autres méthodes.

T. Tris quadratiques.

(1) Tri par insertion

Le tri par insertion est, parmi les algorithmes de tri, un des premiers auquel on pense. On s'en sert notamment pour trier nos jeux de cartes, entre autres.

Il se présente comme suit

Def 1: Le tri par insertion

TRI-INSERTION(A)

pour j entre 1 et longueur(A)

 x = A[j]

 i = j - 1

tant que i > 0 et A[i] > x

 A[i+1] = A[i]

 i = i - 1

A[i+1] = x.

Retourner A.

Pte 8: Le tri par insertion est un tri en place, et un tri stable.

Pte 3

Le tri par insertion possède une complexité en moyenne, et une complexité dans le pire cas, en $O(n^2)$. Dans le meilleur cas, la complexité est en $O(n)$.

(La complexité est ici calculée en comptant les comparaisons dans l'algorithme.)

On a donc une complexité en $O(n^2)$, ce qui est moins intéressant que les algorithmes que nous venons d'écouter.

Toutefois, le tri par insertion se révèle plus efficace, sur de petites entrées, ou des entrées déjà presque triées.

que les algorithmes basés sur la méthode "diviser pour régner."

C'est pour cela qu'il est parfois utilisé en combinaison avec d'autres algorithmes de tri.

② Autres exemples de tris quadratiques.

Le tri par insertion n'est pas le seul tri quadratique, on peut également citer :

- le tri par sélection, dans lequel on recherche à la fois le plus petit élément d'une sous-tableau de l'entrée, pour le mettre au début
- le tri bulle, dans lequel on fait remonter les plus grands éléments d'une liste, en les comparant avec leurs successeurs.

Ces algorithmes ont eux aussi une complexité en $O(n^2)$, contrairement aux tris par comparaison que nous allons étudier désormais.

II Tris dichotomiques

Ces méthodes de tri sont issues de la méthode "diviser pour régner", et s'ils sont aussi des tris par comparaison, ils présentent de meilleures complexités que les tris quadratiques.

① Le tri fusion

Le tri fusion se base sur 2 algorithmes pour marcher.

Def 4:

Tri-fusion (A, p, r)
si $p < r$
$q = \lfloor (p+r)/2 \rfloor$
Tri-fusion (A, p, q)
Tri-fusion (A, q, r)
Fusion (A, p, q, r)

L'algorithme fusion, lui, fusionne deux listes, en cherchant, à chaque itération, le plus petit élément dans les 2 listes, et en l'ajoutant à une nouvelle liste.

Prés 5: Le tri fusion est un tri stable, mais pas un tri en place.

En effet, dans l'algorithme fusion, on recrée une liste pour en fusionner 2.

Prés 6: Le tri fusion possède une complexité en $O(n \log n)$ en moyenne.

Cette complexité est bien meilleure que celle des tris quadratiques, mais peut-on, avec des tris par comparaison, faire mieux que cela?

Avantages du tri fusion: plus rapide, stable

Inconvénient: Algorithme qui demande beaucoup d'espace, non en place.

② Autres tris dichotomiques

Le tri par tas est un autre tri de complexité $O(n \log n)$, qui a l'avantage d'être un tri en place.

Il consiste à voir le tableau en entrée comme un arbre binaire avec le premier élément en racine

consiste à faire remonter le plus gros élément du tableau à la racine, puis ensuite le retirer puis le remplacer par un autre élément, et ainsi de suite

Si il reste un tri en $O(m \log m)$, il reste malgré tout plus lent qu'un autre tri en $O(m \log m)$.

Le tri rapide, qui consiste à placer un pivot à sa bonne position, puis à permuter tous les éléments autour de lui, suivant leur place relative.

Théorème 7: Le tri rapide a une complexité au pire en $O(m^2)$, et une complexité moyenne en $O(n \log m)$ [Duppt 1].

Thm 8: La meilleure complexité que l'on puisse avoir, pour un tri par comparaison, est $O(m \log m)$ [Duppt 2].

Sachant donc cette limite atteinte, existe-t-il d'autres méthodes de tri qui offrent une meilleure complexité.

III: D'autres méthodes de tri

① Le tri par dénombrement

On suppose ici que les éléments de l'entrée sont des entiers de l'intervalle $[0, k]$, avec $k \in \mathbb{N}$.

Si $k = O(m)$, l'algorithme s'exécute en temps $\Theta(m)$.

Le principe du tri par dénombrement consiste à déterminer, pour chaque élément x de l'entrée, le nombre d'éléments inférieurs à x . On peut alors directement placer x au bon endroit.

Tri. Dénombrement s'exécute en 3 temps, pour i entier:

- ① Remplir un tableau avec, en i -ème position, le nombre d'éléments égaux à i .
- ② Remplacer les valeurs de ce tableau par le nombre d'éléments inférieurs ou égaux à i .
- ③ Placer l'élément i à sa position, exacte dans le tableau final.

On a donc un algorithme ayant une complexité en $\Theta(m)$, et qui demande un espace mémoire en $O(m)$.

② Autres tris

Parmi les autres tri dont nous aurions pu parler, on pourra citer

- le tri par bande
- le tri par paquet
- les nouveaux de tri [Cor; ancienne version] (hors-programme).

Reference: Cormen-Leiserson-Rivest-Stein: Algorithmique, 2^{ème} et 3^{ème} éditions [Cormen]
Friedenau-Gaucher-Soria. Types de données et algorithmes: [FGS]

Amere: Algorithmes du tri rapide

TRI-RAPIDE(A, p, r)

$\Delta i \leftarrow p, r$

$q \leftarrow \text{PARTITION}(A, p, r)$

TRI-RAPIDE($A, p, q-1$)

TRI-RAPIDE($A, q+1, r$)

PARTITION(A, p, r)

$x \leftarrow A[p]$

$i \leftarrow p-1$

pour j varie p de $r-1$

$\Delta i \leftarrow A[j] \leq x$

$i \leftarrow i+1$

permuter $A[i]$ et $A[j]$.

permuter $A[i+1]$ et $A[p]$

Retourner $i+1$

Complexité du Quicksort.

(I) Complexité dans le pire des cas.

Étudions tout d'abord le cas du partitionnement.

Le cas le plus défavorable devrait de nature un seul problème à $(n-1)$ éléments.

On suppose que ce cas de figure arrive tout le temps, on a alors, sachant partitionnement coûte $\Theta(n)$

$T(n) = T(n-1) + \Theta(n)$, avec $T(n)$ le nombre de comparaisons effectuées.

Ce qui donne $T(n) = \Theta(n^2)$

Plus mesurement, ~~le cas de figure arrive tout le temps~~

Rec : $T(n) = \max_{0 \leq q \leq n-1} (T(q) + T(n-1-q)) + \Theta(n)$

ou $T(n) \leq cn^2$, pour un certain c constant, d'où

$$T(n) \leq \max_{0 \leq q \leq n-1} (cq^2 + c(n-q-1)^2) + \Theta(n)$$

$$T(n) \leq C \max_{0 \leq q \leq n-1} (q^2 + (n-q-1)^2) + \Theta(n)$$

l'expression $q^2 (n-q-1)^2$ atteint ses extrêmes aux extrêmes

$$\text{d'où } q^2 (n-q-1)^2 \leq (n-1)^2 = n^2 - 2n + 1.$$

$$\text{d'où } T(n) \leq cn^2 - 2cn + 1 + \Theta(n)$$

$$T(n) \leq cn^2$$

on peut choisir c suffisamment grand pour que $\Theta(n)$

$$\text{donc } T(n) = \mathcal{O}(n^2)$$

(II) Complexité moyenne

Pour simplifier les notations, on appelle S_1, \dots, S_m les éléments du tableau, et Z_q

Mais allons donc compter le nombre de comparaisons effectuées dans l'exécution

$$X = \sum_{i=1}^{n-1} \sum_{j=i+1}^n X_{ij}$$

et donc $E[X] = \sum_{i=1}^{n-1} \sum_{j=i+1}^n E[X_{ij}]$ par linéarité de l'espérance

$$E[X] = \sum_{i=1}^{n-1} \sum_{j=i+1}^n \mathbb{P}\{z_i \text{ et } z_j \text{ sont comparés, dans } Z_{ij}\}.$$

Or z_i et z_j sont comparés, dans Z_{ij} , uniquement si l'un de z est le premier pivot choisi, sinon, ils seront séparés dans z dans tableaux et ne seront plus comparés par la suite.

Donc $\mathbb{P}\{z_i \text{ et } z_j \text{ sont comparés, dans } Z_{ij}\} = \mathbb{P}\{z_i \text{ est le premier pivot} \mid z_i \text{ est le } i^{\text{ème}}$

des événements étant indépendants; $\mathbb{P}\{z_i \text{ et } z_j \text{ sont comparés}\} = \mathbb{P}\{z_i = 1^{\text{ère}} \mid i^{\text{ème}}\} + \mathbb{P}\{z_j = 1^{\text{ère}} \mid j^{\text{ème}}\}.$

Les pivots étant choisis aléatoirement, avec distribution uniforme, ces probabilités valent $\frac{1}{j+1-i}$

$$\text{donc } \mathbb{P}\{z_i \text{ et } z_j \text{ sont comparés}\} = \frac{2}{j-i+1}$$

$$\text{donc } E[X] = \sum_{i=1}^{n-1} \sum_{j=i+1}^n \frac{2}{j-i+1}$$

Par changement de variable, on a. ($k=j-i+1$).

$$E[X] = \sum_{i=1}^{n-1} \sum_{j=i+1}^n \frac{2}{j-i+1} = \sum_{i=1}^{n-1} \sum_{k=2}^{n-i+1} \frac{2}{k} < \sum_{i=1}^{n-1} \sum_{k=2}^n \frac{2}{k}$$

$$E[X] < \sum_{i=1}^{n-1} O(n \log n) = O(n \log n)$$

On arrive bien au temps attendu de $O(n \log n)$.

Recherche: [Complexité]

Dans le pire cas? Sans la complexité quadratique

↳ à chaque fois, on choisit le pire des voisinages ou le minimum.

$$\mathbb{P}(c = O(n^2)) = \frac{2}{n} \mathbb{P}(c = O((n-1)^2)) = \frac{2^{n-1}}{n}$$

Optimalité du tri par comparaison

Rappel: La meilleure complexité, pour un tri par comparaison, est $O(n \log n)$.

On parle ici de comparaison au pire et en moyenne, et on compte les comparaisons. Bien entendu, celle bonne profondeur, on va utiliser les arbres binaires de décisions.

↳ Lemme 1: L'arbre binaire de décision associé à un algo de tri par comparaison a exactement $n!$ feuilles.

↳ Preuve: Chaque feuille indique le résultat des différentes exécutions de toutes les données possibles.

Le nombre d'ordres possibles pour les n éléments dont $n!$, on a le proposé.

Autres $TC(n)$ le nombre minimal de comparaison, (et $TC_{\text{moy}}(n)$ le nombre moyen) a effectivement dans le pire cas (resp. en moyenne)

on a alors $TC_{\text{max}}(n) = \inf \{R(A), A \text{ arbre de décision associé à l'algo de tri}\}$

$TC_{\text{moy}}(n) = \inf \{P(A), A \text{ arbre de décision associé à l'algo de tri, avec } P(A) \text{ profondeur moyenne externe de } A\}$.

↳ Lemme 2: Tout arbre binaire de décision ayant n feuilles a une hauteur égale à $\lceil \log_2 n \rceil$

On en déduit que $TC_{\text{max}}(n) \geq \lceil \log_2(n) \rceil$

De plus $n! \leq n^n$

on a $T_{\log(m)} = \Theta(m \log m)$.

Et on connaît des algos de tri dont la complexité au pire est $\Theta(m \log m)$.

Exemple le tri fusion pour lequel le nombre de comparaisons par bien n éléments réelle

la relation f sachant que dans le pire des cas, il y a $n-1$ comparaisons λ

$$F(n) = n-1 + F(\lfloor n/2 \rfloor) + F(\lceil n/2 \rceil)$$

et alors $F(n) = n \log_2(n) - 2^{\lceil \log_2 n \rceil} + 1$.

Pour le cas de la complexité moyenne

on a $T_{\text{moy}}(n) = \inf \{ \{ \Theta(n) \} \}$.

Lemme 3: Pour un arbre binaire A ayant n feuilles, $\rho(A) \geq \log_2(n)$.

Ce lemme nous donne alors directement $T_{\text{moy}}(n) \geq \log_2(n)$, et donc

$$T_{\text{moy}}(n) = \Theta(m \log m)$$

Un Pochonelle, on obtient le tri rapide, dont on sait qu'il atteint cette borne

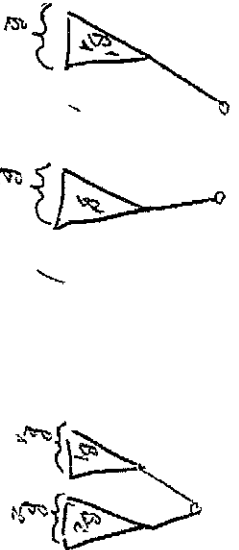
/// Rem du Lemme 3 : raisonnablement par récurrence.

Inductif: c'est vrai pour l'arbre à une feuille qui est de profondeur 0.

Supposons la propriété vraie pour les arbres ayant moins de $k-1$ feuilles

Soit B un arbre à k feuilles

Est d'une de ces 3 formes



Le seul cas a montrer ici est le 3^e cas.

Comme $B_1 < B$ et $B_2 < B$, q_1 et q_2 vérifient la propriété par récurrence, donc

$$PE(q_1) \geq \log_{B_1} B_1 \quad \text{et} \quad PE(q_2) \geq \log_{B_2} B_2.$$

$$\text{or } PE(B) = 1 + \frac{B_1}{B_1 + B_2} PE(q_1) + \frac{B_2}{B_1 + B_2} PE(q_2).$$

$$\text{donc } PE(B) = 1 + \frac{B_1 \log_{B_1} B_1 + B_2 \log_{B_2} B_2}{B_1 + B_2}.$$

Comme $B_1 + B_2 = B$, on peut écrire

$$PE(B) = 1 + \frac{B_1 \log_{B_1} B_1 + (B - B_1) \log_{B - B_1} (B - B_1)}{B} = q(B_1).$$

et $q(B_1)$ est maximal quand $B_1 = \frac{B}{2}$, et alors $q(B/2) \geq \log_{B/2} (B/2)$

ce qui achève la récurrence.

Référence: [FGS].