

Complexité algorithmique

Florent Bouchez Tichadou

5 septembre 2019

L'algorithmique est la science qui s'intéresse non seulement à l'écriture des algorithmes, mais également à leur étude et analyse. Dans ce document, nous abordons la notion de *complexité algorithmique*, qui est une mesure de l'« efficacité » d'un algorithme. Nous nous intéressons donc non seulement à l'écriture d'algorithmes qui produisent des résultats corrects, mais également à la vitesse à laquelle ils résolvent le problème.

Note importante : contrairement à ce que le nom suggère, la complexité *n'est pas* une mesure de si un algorithme est « simple » ou « complexe » d'un point de vue humain. C'est en fait bien souvent l'inverse : un algorithme simple aura généralement une complexité plus élevée (il « prend plus de temps ») qu'un algorithme ingénieux, qui aura une faible complexité (plus « rapide »).

1 Introduction à la complexité algorithmique

La *complexité d'un algorithme* est une prédiction ou une garantie que l'algorithme ne prendra jamais plus qu'un certain nombre d'étapes ou opérations, qui dépend souvent de la *taille* des données qu'il manipule. On note en général n cette taille et on cherche « formule(n) » qui représente le nombre maximum d'opérations et dépend de l'algorithme.

Voyons un algorithme simple qui renvoie la somme de tous éléments d'un ensemble représenté dans un tableau avec longueur explicite.

```
-----  
-----  
somme(E)  
┌ s ← 0  
├ pour i de 0 à E.longueur - 1 faire  
│   ┌ s ← s + E.tab[i]  
└   └ retourner s  
-----  
-----
```

Analysons les différentes parties de cet algorithme :

- l'initialisation de s coûte 1 opération ;
- pour la boucle, le nombre d'opérations est la somme des opérations de chacune des itérations :
 - corps de boucle : 2 opérations : une addition puis le stockage du résultat dans s ;
 - maintenance de boucle : à chaque fois un incrément et un test sont cachés dans le « pour » :
 $i \leftarrow i + 1$ et $i \leq \text{longueur} - 1$: 2 opérations ;
 - entrée dans la boucle : le premier test (1 opération, pas d'incrément au début) ;
- retour de la fonction : 1 opération.

Au final, chaque itération de la boucle exécute le même nombre d'instructions, le nombre total pour la boucle est donc « coût d'une itération » \times « nombre d'itérations » ; la boucle est répétée « longueur » fois. La formule pour cet algorithme est $1 + 1 + n \times (2 + 2) + 1 = 3 + 4n$ où $n = \text{longueur de l'ensemble}$.

Nous voyons ici que le temps d'exécution de cet algorithme croît linéairement en n . Quand nous étudions la complexité algorithmique, ce qui nous intéresse le plus est le comportement *asymptotique* de l'algorithme, c'est-à-dire la « forme » de la formule pour n grand. Nous gardons donc seulement le « plus gros » terme de la formule sans sa constante : nous utilisons la notation « grand O » des mathématiques, et disons que l'algorithme est en $O(n)$.

2 La notation grand O

Étant donnés deux fonctions f et g , $f = O(g)$ ssi pour tout n suffisamment grand, il existe une constante C telle que $f(n) \leq C \times g(n)$. Plus formellement :

$$\exists n_0 \geq 0, \exists C > 0 \text{ t.q. } \forall n \geq n_0, f(n) \leq C \times g(n)$$

Cette notation nous permet de nous concentrer sur les « grandes » valeurs de n , sans se préoccuper des petites valeurs de n , pour lesquelles de toutes façons l'algorithme se terminera très rapidement sur un ordinateur moderne.

La constante C de la notation nous permet d'« oublier » la constante du plus grand facteur. L'analyse devient non seulement plus simple mais surtout possible car il est en pratique presque impossible de savoir exactement combien d'opérations sont nécessaires : cela dépend notamment du processeur où le code est exécuté. Par exemple dans le corps d'une boucle, ce qui nous importe est alors juste de savoir qu'il y a un nombre constant d'opérations, et si la boucle est répétée n fois on a alors $O(n)$ opérations.

Ces deux propriétés combinées nous permettent de nous concentrer sur le comportement pour des « grandes » instances, où il y a une énorme différence par exemple entre une complexité en $O(n)$ et une en $O(n^2)$, mais où des complexités qui seraient $\approx 4 \times n$ et $\approx 5 \times n$ sont considérées comme équivalentes.

Revenons à notre exemple pour en faire directement l'analyse avec la notation O . Pour mémoire, rappelons qu'un nombre constant d'opération est en $O(1)$. Nous avons maintenant les propriétés suivantes pour une analyse de la complexité en temps. Le détail est à gauche mais on général on annote directement sur l'algorithme comme montré ci-dessous à droite.

- initialisation : $O(1)$;
- boucle : n fois le corps de boucle ;
- corps de boucle : $O(1)$;
- finalisation : $O(1)$.

somme(E)	
s ← 0	$O(1)$
pour i de 0 à E.longueur -1 faire	$n \times$
└ s ← s + E.tab[i]	$O(1)$
retourner s	$O(1)$

Au global, la complexité est $O(1) + n \times O(1) + O(1) = O(1) + O(n) = O(n)$.

3 Comparer les complexités

Le but de l'analyse de complexité est de pouvoir comparer plus facilement différents algorithmes qui effectuent la même tâche. On préférera par exemple l'algorithme qui s'exécute en $O(n)$ par rapport à celui en $O(n^2)$. Il devient également possible de prédire le temps que prendra un algorithme à trouver la solution sur une instance très grande en extrapolant une mesure de temps d'exécution faite sur une plus petite instance.

Dans cette section, nous nous intéressons à observer combien de fois plus de temps est nécessaire à un algorithme quand on double la taille des données d'entrée. Par exemple, un algorithme linéaire ($O(n)$) mettra deux fois plus de temps. On dénote $\log n$ le logarithme en base 2, qui fait partie des complexités courantes pour un algorithme (voir sections suivantes). Nous donnons ici les complexités dans l'ordre du plus courant au moins courant (sans être exhaustif, il existe bien sûr beaucoup d'autres complexités possibles). Nous donnons également une idée de la taille maximale que peuvent avoir les données (colonne « Max n ») avant que l'algorithme devienne impraticable (cela prendrait trop de temps pour résoudre le problème).

Complexité	Nom courant	Temps quand on double la taille de l'entrée	Max n
$O(n)$	linéaire	prend 2 fois plus de temps	10^{12}
$O(1)$	constant	prend le même temps	pas de limite
$O(n^2)$	quadratique	prend 4 fois plus de temps	10^6
$O(n^3)$	cubique	prend 8 fois plus de temps	10 000
$O(\log n)$	logarithmique	prend seulement une étape de plus	$10^{10^{12}}$
$O(n \log n)$	linearithmique	prend deux fois plus de temps + $\log n$	10^{11}
$O(2^n)$	exponentiel	prend tellement de temps que c'est inconcevable	30

On observe une séparation nette entre la complexité exponentielle et les autres, qu'on appelle *polynomiales*. Les algorithmes exponentiels sont si catastrophiques qu'il ne sert à rien de se demander « que se passe-t-il si je double la taille de mon entrée » car, à l'inverse, il suffit d'augmenter la taille du 1 pour doubler le temps de calcul ! Ces algorithmes doivent être évités à tout prix, sauf si l'on sait que les instances sont très petites.

L'algorithme utilise deux boucles imbriquées. La boucle extérieure a n itération, et celle intérieure v itération, avec v qui varie : $n - 1$ au début mais diminue ensuite au cours du temps jusque 1. Nous savons que $v \leq n$ est toujours vrai, donc nous savons que le corps de la boucle interne est exécuté à $O(n)$ fois à chaque itération. Le corp boucle est en $O(1)$ (ainsi que le reste des calculs), donc la complexité globale est $O(n^2)$.

De même que vu précédemment, on pourrait faire une analyse plus précise qui nous conduirait à trouver qu'on exécute le corps de la boucle interne $n \times (n - 1)/2$ fois, mais cela ne changerait pas la complexité.

5.3.2 Dans une liste chaînée

```

tri_selection(S)
  dernier ← Nil
  tant que dernier ≠ S.tête faire
    max ← S.tête.valeur
    maxcel ← S.tête
    maxpred ← Nil
    cel ← maxcel.suivant
    queue ← maxcel
    tant que cel ≠ dernier faire
      si cel.valeur > max alors
        max ← cel.valeur
        maxcel ← cel
        maxpred ← queue
      queue ← cel
      cel ← cel.suivant
    si queue ≠ maxcel alors
      si maxpred = Nil alors
        S.tête ← maxcel.suivant
      sinon
        maxpred.suivant ← maxcel.suivant
    maxcel.suivant ← dernier
    queue.suivant ← maxcel

```

Complexity analysis for the code above:

- $O(1)$ for `dernier ← Nil`
- $n \times$ for the outer loop `tant que dernier ≠ S.tête faire`
- $O(1)$ for `max ← S.tête.valeur`
- $O(1)$ for `maxcel ← S.tête`
- $O(1)$ for `maxpred ← Nil`
- $O(1)$ for `cel ← maxcel.suivant`
- $O(1)$ for `queue ← maxcel`
- $O(n) \times$ for the inner loop `tant que cel ≠ dernier faire`
- $O(1)$ for `si cel.valeur > max alors`
- $O(1)$ for `max ← cel.valeur`
- $O(1)$ for `maxcel ← cel`
- $O(1)$ for `maxpred ← queue`
- $O(1)$ for `queue ← cel`
- $O(1)$ for `cel ← cel.suivant`
- $O(1)$ for `si queue ≠ maxcel alors`
- $O(1)$ for `si maxpred = Nil alors`
- $O(1)$ for `S.tête ← maxcel.suivant`
- $O(1)$ for `sinon`
- $O(1)$ for `maxpred.suivant ← maxcel.suivant`
- $O(1)$ for `maxcel.suivant ← dernier`
- $O(1)$ for `queue.suivant ← maxcel`

Analyse de complexité : chaque boucle « tant que » sur la liste chaînée agit de manière similaire à la boucle « pour » correspondante dans la représentation par tableau. Comme précédemment, on ne connaît pas exactement le nombre d'itérations de la boucle interne, mais celui-ci est borné par n donc $O(n)$ itérations, et au final la complexité est $O(n^2)$.

5.4 Recherche d'un élément dans une séquence

Considérons deux cas : selon que la séquence est triée ou non.

Séquence non triée Il nous faut vérifier chaque élément de la séquence : c'est un parcours classique, donc $O(n)$ pour un tableau ou une liste chaînée.

Séquence triée Pour une liste chaînée, cela nous permet de nous arrêter dès que l'on trouve un élément strictement plus grand que celui que l'on cherche. Cela ne change tout de même pas la complexité en général, car dans le pire des cas, il nous faut toujours vérifier jusqu'au dernier élément de la séquence. En moyenne on vérifiera la moitié de la séquence, ce qui est toujours du $O(n)$.

Pour un tableau, nous pouvons utiliser la *dichotomie* puisque nous avons directement accès au « milieu » de la séquence : en comparant avec l'élément du milieu, on n'a plus ensuite qu'à chercher dans une moitié de la séquence. Supposons que l'on cherche x , et que $S.tab[\text{longueur}/2]$ contient v ; si $x < v$, alors x ne peut pas être à un indice plus grand que $\text{longueur}/2$. En utilisant cette propriété, on peut alors faire une *recherche binaire*, algorithme présenté à droite.

```

recherche (S, x)                               /* S doit être trié */
┌
│  i ← 0                                       O(1)
│  j ← S.longueur - 1                         O(1)
│  tant que i ≤ j faire                       k ×
│  │  m ← (i + j) div 2                       O(1)
│  │  si x = S.tab[m] alors retourner m       O(1)
│  │  sinon si x < S.tab[m] alors j ← m - 1   O(1)
│  │  sinon i ← m + 1                         O(1)
│  retourner -1                               /* non trouvé */
└

```

Quelle est la complexité de cet algorithme? Nous voyons que les parties avant et après la boucle sont en $O(1)$, de même que le corps de boucle. La complexité va donc être directement liée au nombre d'itérations de cette boucle que l'on note k .

Considérons la quantité $j - i$: au début, elle est égale à $S.\text{longueur}$, et à chaque itération cette quantité va être divisée approximativement par 2. Dans le pire des cas (quand x n'est pas dans la séquence), ce processus est répété jusqu'à atteindre un, puis zéro, puis i devient plus grand que j (à cause du -1 ou du $+1$).

Nous devons donc répondre à la question suivante : combien de fois faut-il diviser n par 2 pour attendre 1 (on peut ne pas considérer les deux dernières étapes car on cherche une notation O)? La réponse est $O(\log n)$ (le logarithme en base 2). En effet, on cherche le plus petit k tel que $n/2^k \leq 1$, c'est-à-dire, $n \leq 2^k$, et nous savons que $2^{\log n} = n$.

Une autre façon de voir (ou prédire) la complexité logarithmique est de remarquer que même si l'on double la taille de l'entrée (i.e., une séquence de taille $2n$ au lieu de n), on n'a besoin que d'une étape de plus pour effectuer la recherche...

En général, on retrouve une complexité logarithmique dans tous les algorithmes qui contiennent une boucle divisant une quantité de donnée par une constante à chaque itération.

5.5 Calcul des nombres de Fibonacci

Les nombres de Fibonacci sont définis ci-dessous à gauche. Un algorithme très naïf pour calculer le $n^{\text{ème}}$ nombre est proposé à droite.

$$\begin{aligned}
 F_0 &= 1 \\
 F_1 &= 1 \\
 F_n &= F_{n-1} + F_{n-2}
 \end{aligned}$$

```

Fibo(n)
┌
│  si n < 2 alors
│  │  retourner 1
│  sinon
│  │  retourner Fibo(n - 1) + Fibo(n - 2)
└

```

Il est plus difficile de calculer la complexité de cet algorithme car il est récursif : la complexité de Fibo dépend... de la complexité de Fibo!

Définissons alors C_n comme étant le nombre d'opérations nécessaire pour calculer $\text{Fibo}(n)$. On peut écrire la formule suivante : $C_n = O(1) + C_{n-1} + C_{n-2}$. Remarquez que cette formule ressemble très fortement à l'équation de récurrence de la suite de Fibonacci elle-même!

Faisons maintenant une approximation. Puisque $\text{Fibo}(n - 1)$ va elle-même appeler $\text{Fibo}(n - 2)$, on sait que $C_{n-1} \geq O(1) + C_{n-2}$, donc on écrit $C_n \leq 2 \times C_{n-1}$. Avec $C_0 = C_1 = 1$, on a directement que $C_n \leq 2^n$, donc la complexité de Fibo est $O(2^n)$.

Bien sûr, cette borne n'est pas exacte puisque nous avons fait une approximation grossière. La complexité réelle pourrait être bien plus petite! Pour prouver que cet algorithme n'est pas polynomial, nous pouvons également dire que $C_n \geq 2 \times C_{n-2}$ et donc $C_n \geq 2^{n/2}$. Cet algorithme a donc une complexité exponentielle, comprise entre $O(2^{n/2})$ et $O(2^n)$.

5.6 Tri fusion

Nous allons voir un algorithme différent pour trier une séquence, basé sur l'idée suivante : si la séquence est de longueur supérieure à 2, on divise la séquence en deux moitié de taille identiques, on trie (récursivement) les deux sous-séquences, puis on les fusionne pour obtenir la séquence triée.

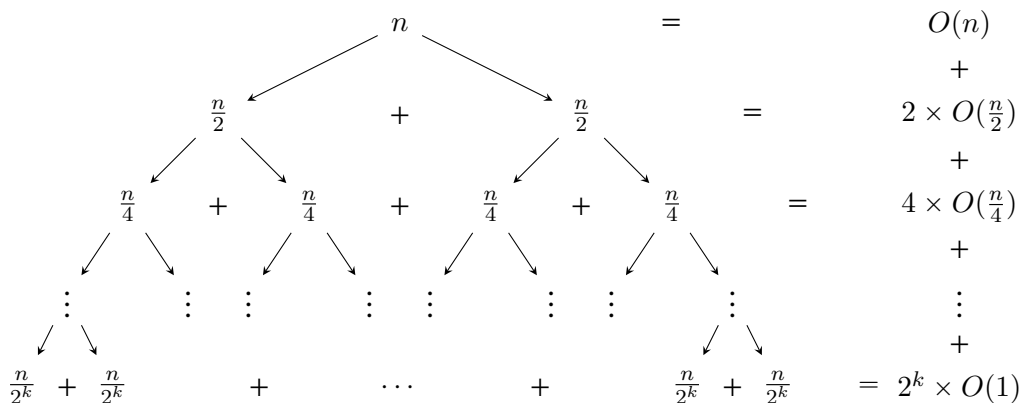
L'implémentation présentée ici n'est pas la plus optimisée, il serait par exemple possible d'effectuer moins de copies entre tableaux. Mais cela rend difficile la lecture de l'algorithme sans pour autant en diminuer la complexité. Nous donnons une implantation à base de tableaux mais il est possible d'utiliser des listes chaînées.

<pre> Tri(S) ┌ si S.longueur = 1 alors │ ┌ retourner │ m ← S.longueur div 2 │ S₁ ← m premiers éléments de S │ S₂ ← les autres │ Tri(S₁) │ Tri(S₂) │ Fusion(S, S₁, S₂) </pre>	<pre> Fusion(S, S₁, S₂) ┌ i ← 0; i₁ ← 0; i₂ ← 0 tant que i₁ < S₁.longueur et i₂ < S₂.longueur faire ┌ si S₁.tab[i₁] < S₂.tab[i₂] alors │ S.tab[i] ← S₁.tab[i₁] │ i₁ ← i₁ + 1 ┌ sinon │ S.tab[i] ← S₂.tab[i₂] │ i₂ ← i₂ + 1 ┌ i ← i + 1 si i₁ = S₁.longueur alors ┌ échanger S₁ et S₂ ┌ échanger i₁ et i₂ pour j de i₁ à S₁.longueur - 1 faire ┌ S.tab[j] ← S₁.tab[j] ┌ i ← i + 1 S.longueur ← S₁.longueur + S₂.longueur </pre>
--	--

L'algorithme de tri fusion utilise une stratégie classique dite « diviser pour régner ». Elle consiste en découper un gros problème en plus petits problèmes, résoudre récursivement les problèmes, puis combiner les résultats pour obtenir une solution au problème initial. Analysons la complexité de cet algorithme.

La fonction Fusion utilise la fonction Tri, nous allons donc d'abord analyser cette dernière. Il y a deux boucles : un « tant que » et un « pour », et on ne sait pas à l'avance combien d'opération chacune d'elle va effectuer. En revanche, on sait que le nombre *total* d'itérations combiné des deux boucles est exactement $n = S_1.\text{longueur} + S_2.\text{longueur}$. Puisque les deux boucles font un nombre constant d'opérations, la complexité de la fonction est $O(n)$. En effet, nous copions ici tous les éléments de S_1 et S_2 vers S exactement une fois.

Voyons maintenant la fonction Tri, qui est récursive. Comme dans la section précédente, écrivons la formule. Si C_n est le nombre d'opérations faites par Tri sur une séquence de taille n , nous avons $C_n = O(n) + 2 \times C_{n/2} + O(n)$. Le premier $O(n)$ correspond à la séparation en S_1 et S_2 et le deuxième leur fusion. Nous allons maintenant dessiner « l'arbre d'appels » d'une exécution de Tri, où la valeur de chaque nœud représente la taille de la séquence d'un appel (récursif) à Tri.



Le coût du tri est la somme des coûts de tous les nœuds de cet arbre d'appels. Le coût d'un nœud est maintenant simplement le coût de la séparation et de la fusion, puisque le coût des deux appels récursifs est maintenant compté dans les nœuds des fils.

Si l'on regarde les coûts par niveau, on se rend compte que le premier niveau est $O(n)$, le deuxième deux fois $O(n/2)$, le troisième quatre fois $O(n/4)$, etc. Chaque niveau coûte donc $O(n)$, où n est la taille de la séquence initiale, jusqu'au dernier niveau, où il y a 2^k feuilles représentant le coût de « trier » des sous-séquences de taille 1, i.e., $O(1)$, et $2^k = n$.

Au total, la complexité de l'algorithme du tri-fusion est $O(n) \times$ le nombre de niveaux, c'est-à-dire la hauteur de l'arbre. Comme pour la recherche binaire, on trouve que cette hauteur est exactement le nombre de fois qu'il faut diviser n par 2 pour atteindre 1, soit $k = \log n$. La complexité du tri fusion est donc en $O(n \log n)$, ce qui est **bien meilleur** que le tri par sélection analysé plus précédemment.