

*Structures de données et complexité*

**Informatique scientifique pour le Calcul**  
**Ecoles Doctorales 2014/2015**

Vincent Miele

CNRS - Biométrie & Biologie Evolutive

Décembre 2014

Document disponible sur <http://lyoncalcul.univ-lyon1.fr>

- ▶ quel(s) langage(s) ? [cours suivant]
- ▶ quelles algorithmes ? adaptés au calcul parallèle ? [cours suivant] comment quantifier les ressources qu'ils demandent ? [ce cours]
- ▶ quelles structures de données pour mon problème (indépendant du langage) ? la transcription de celles que j'utilise sur la papier ? [ce cours]



## Objectifs

- ▶ efficacité → bien comprendre les coûts de calcul, théoriques [ce cours] et coût machine [cours précédent Violaine]
- ▶ robustesse → utiliser au maximum les standards, les bibliothèques, s'appuyer sur la communauté d'utilisateurs [cours suivant]
- ▶ maximum de cohérence avec les algorithmes du papier → lisibilité, encapsulation [ce/autre cours]

### Définition

La *complexité* d'un problème caractérise les ressources nécessaires pour les résoudre. Les ressources essentielles sont le temps (d'exécution) et l'espace (mémoire).

Jusqu'aux années 70, seule la mesure expérimentale de la complexité d'un algorithme était (parfois) effectuée et dépendait des machines. . . une mesure théorique devint souhaitable !

### Définition

La *complexité algorithmique* est un ordre de grandeur théorique du temps de calcul (*complexité en temps*) et/ou de l'espace mémoire (*complexité en mémoire*) utilisé en fonction d'une mesure des données.

Souvent, si un algorithme permet de gagner du temps de calcul, il occupe davantage de place en mémoire. . . vigilance. . .

La complexité n'est pas la même selon les cas

- ▶ complexité au pire : complexité maximum, dans le cas le plus défavorable
- ▶ complexité en moyenne : il s'agit de la moyenne des complexités obtenues selon les cas
- ▶ complexité au mieux : complexité minimum, dans le cas le plus favorable.  
En pratique, cette complexité n'est pas très utile

Exemple (complexité en temps) : chercher un élément dans une liste de taille  $n$

pire :  $n$  opérations

moyenne :  $(n+1)/2$  opérations

au mieux : 1 opération

## └ Complexité

### └ Notation de Landau

Si on veut comparer des algorithmes sans considérer l'implémentation, on peut comparer les complexités au regard de la taille attendue des données

- ▶ complexité au pire
- ▶ forme générale de la complexité, qui indique la façon dont elle évolue en fonction d'un paramètre (ou plusieurs).

On utilisera la notation  $O(\dots)$  qui veut dire *de l'ordre de...*

#### Definition

*Notation de Landau* :  $g$  est dite en  $O(f)$  s'il existe des constantes  $c > 0$  et  $x_0 \gg 0$  tq  $g(x) < c * f(x)$  pour tout  $x > x_0$ . On note  $g = O(f)$  et on dit que  $g$  est dominée asymptotiquement par  $f$ .

Si  $g = O(f)$  et  $f = O(h)$  alors  $g = O(h)$

si  $g = O(f)$  et  $k$  un nombre, alors  $k * g = O(f)$

si  $f_1 = O(g_1)$  et  $f_2 = O(g_2)$  alors  $f_1 + f_2 = O(g_1 + g_2)$

si  $f_1 = O(g_1)$  et  $f_2 = O(g_2)$  alors  $f_1 * f_2 = O(g_1 * g_2)$

Si  $n$  est la taille des données, complexité en temps :

$O(1)$  *constante* : extraction d'un élément dans un tableau

$O(\log(n))$  *logarithmique* : recherche dichotomique dans une liste triée (voir après)

$O(n)$  *linéaire* : parcours d'une liste

$O(n \log(n))$  *quasi-linéaire* : tri d'un tableau

$O(n^2)$  *quadratique* : calcul du maximum d'une matrice carré

$O(n^3)$  *cubique* : multiplication matricielle (voir après)

$O(2^n)$  *exponentielle* : problème du sac à dos par force brute (énumération)...

NB : la notation de Landau peut utiliser plusieurs paramètres décrivant la taille du problème (exple :  $O(m * n + q)$ )

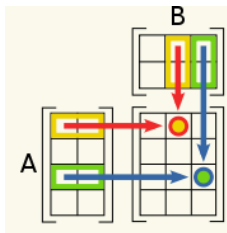
## Complexité

### A l'attaque!

## Produit matrice x matrice

Pour la produit de deux matrice de taille  $m \times n$  et  $n \times p$ , on réalise  $n$  produits  $+n$  sommes pour chaque élément de la matrice finale de taille  $m \times p$ .

On a donc une complexité  $O(m * p)$  en espace et  $O(n * m * p)$  en temps



Pour le cas où  $m = p = n$ ,  $O(n^2)$  en espace et  $O(n^3)$  en temps.

“L’algorithme de Coppersmith-Winograd [...] est en  $O(n^{2,376})$  [...] Mais aucune implémentation de l’algorithme n’est utilisée car la constante dans le grand O est prohibitive.” (Wikipedia)



## └ Complexité

## └ A l'attaque!

## Recherche dichotomique dans un tableau trié

On coupe le tableau en deux et on cherche l'élément dans une des deux parties en répétant le même traitement récursivement.

```
fonction rechercheElement(tab, x)
  i <- 0; j <- tab.longueur-1;
  tantque (i <= j) faire
    si (tab[(j+i)/2] = x) alors retourne VRAI;
    sinon
      si (tab[(j+i)/2] > x) alors j <- (j+i)/2 - 1;
      sinon i <- (j+i)/2 + 1;
```

Au pire, la longueur de la partie du segment de tableau  $[i, j]$  à traiter est d'abord  $n$ , puis  $n/2$ , puis  $n/4$ , ..., jusqu'à ce que  $n/2^t = 1$ .

Le nombre de tours de boucles est  $t = \log_2(n)$ . La complexité en temps est donc  $O(\log(n))$ .  $O(n)$  en espace.

Le nombre de processeurs disponibles est une donnée à prendre en compte dans la complexité. Deux alternatives pour exprimer la complexité :

- ▶  $O(f(n))$  avec  $g(n)$  processeurs : exprime la complexité obtenue avec un nombre idéal de processeurs
- ▶  $O(f(n, p))$  avec  $p$  le nombre de processeurs : plus pragmatique  
Permet d'intégrer la loi d'Amdhal :  $O_{seq}(f(n)) + O_{par}(g(n)/p)$   
Si  $O_{seq}(f(n))$  est dominé, alors on garde  $O_{par}(g(n)/p)$ .

NB : attention à  $O_{comm}(h(n))$  !

└ Complexité

└ La complexité comme indicateur

---

La complexité nous aide à choisir le bon algorithme et les bonnes structures de données. C'est un indicateur.

Mais attention aux constantes devant le  $O(\dots)$ !

Et attention aux coûts réels engendrés par les accès mémoire, la bande passante, les architectures des processeurs.

Un bon algorithme en machine est d'abord celui qui privilégie la localité spatiale et temporelle !

## └ Complexité

## └ La complexité comme indicateur

Deux algorithmes de même complexité, mais l'un respecte la localité spatiale et temporelle, l'autre non.

```
1 #include <vector>
2 using namespace std;
3 int main(int argc, char ** argv )
4 {
5     int s = 18000;
6     vector<vector<double>> v(s, vector<double>(s, 0.));
7
8     // real 0m8.970s
9     for (int i=0; i<s; i++)
10         for (int j=0; j<s; j++)
11             v[i][j] = v[i][j] + v[i][j];
12
13     // real 0m45.077s
14     for (int i=0; i<s; i++)
15         for (int j=0; j<s; j++)
16             v[j][i] = v[j][i] + v[j][i];
17 }
```

## └ Structures de données

### └ Du cerveau à l'ordinateur : copier/coller ?

---

Prélude : dans ce qui suit, on s'appuie sur les *containeurs* de la STL C++, mais des équivalents existent dans les autres langages.

Sur la papier, on a des structures de données de type :

- ▶ vecteur
- ▶ matrice pleine
- ▶ matrice creuse (ou graphe) avec  $m$  éléments tq  $m \sim O(n)$  ou  $m \ll n^2$
- ▶ des listes (triés), des piles
- ▶ des couples “identifiants-observations”
- ▶ ...

Doit-on chercher à utiliser formellement les mêmes structures sur le papier et en machine ? ou seulement moralement ?

On veut des performances en temps ET de la parcimonie en espace.

La clé de voûte de l'efficacité, car cohérence avec la structure en lignes de caches (cf. cours précédent)

- ▶ simple tableau de taille  $n$  (mémoire contigüe!) avec information de longueur
- ▶ access  $O(1)$
- ▶ insert  $O(n)$ ,  $O(1)$  à la fin
- ▶ find  $O(n)$
- ▶ delete  $O(n)$ ,  $O(1)$  à la fin
- ▶ sort  $O(n \log(n))$

NB : on peut aussi créer des tableaux comme en C

## Encapsulation d'un vecteur interne pour gérer les matrices pleines

```
1 class Matrix
2 {
3     private:
4         int _s;
5         vector<double> _vinternal;
6     public:
7         Matrix(int s)
8             : _vinternal(s*s, 0), _s(s) {}
9         Matrix(const Matrix& m)
10            : _vinternal(m._vinternal), _s(m._s) {}
11         ~Matrix() {}
12         inline double& element(int i, int j){
13             return _vinternal[i*_s+j];
14         }
15 };
```

Structure favorisant l'ajout/retrait des éléments. La notion de longueur n'est pas fondamentale car variable.

- ▶ liste doublement-chainée
- ▶ access  $NA$
- ▶ insert  $O(1)$
- ▶ find  $O(n)$
- ▶ delete  $O(1)$
- ▶ sort  $O(n\log(n))$



**Exple :** Une liste d'éléments évolue, chacun pouvant survivre ou non (delete de liste) et donner naissance à un nouvel élément (ajout en fin)

- ▶ une `list<.>` devrait convenir
- ▶ une approche basée sur les `vector<.>` mérite d'être testée, avec ajout en fin de vecteur (naissance) ou copie du dernier élément à la place du delete (non survie) → danger de la copie

A tester après profiling !

## └ Structures de données

└ Liste triée - `set<.>`

Permet de gérer un ensemble (au sens ensembliste - unicité) d'éléments qui seront placés dans la structure selon un critère de tri.

Premier conteneur associatif, i.e. on accède à l'élément par une clé.

- ▶ red-black tree (variante des arbres binaires<sup>1</sup>) [▶ Démo](#)
- ▶ access *NA*
- ▶ insert  $O(\log(n))$
- ▶ find  $O(\log(n))$
- ▶ delete  $O(\log(N))$
- ▶ sort *NA*

Structure triée par construction, mais as-t-on besoin qu'elle le soit à chaque étape de sa construction ?

---

1. plus petites (grandes) valeurs dans le sous arbre de gauche (droite resp.)

## └ Structures de données

## └ Liste triée - set&lt;.&gt;

Exple : Insertion dans un set<.> ou tri/unicité à la fin avec un vector<.>?

```
1 #include <set>
2 #include <vector>
3 #include <algorithm>
4 #include <stdlib.h>
5 using namespace std;
6 int tirage(int n){
7     return (int)((double)rand()/((double)RAND.MAX + 1) * n);
8 }
9 int main(int argc, char ** argv)
10 {
11     vector<int> vec;
12     set<int> s;
13     int n=1e7;
14
15     // 2 sec
16     for (int i=0; i<n; i++)
17         vec.push_back(tirage(n));
18     sort(vec.begin(), vec.end());
19     vec.erase( unique( vec.begin(), vec.end() ), vec.end() );
20
21     // 16 secondes
22     for (int i=0; i<n; i++)
23         s.insert(tirage(n));
24 }
```

## └ Structures de données

## └ Liste triée - set&lt;.&gt;

**Exple :** Une matrice creuse en 0/1 (cas d'école ici) est une liste d'éléments non nuls avec leur index. Comment l'interroger ? A-t-on besoin de la modifier ?

- ▶ `list<pair<unsigned int, unsigned int> >` : liste des paires d'indices  
→ difficile à interroger
- ▶ `list<list<unsigned int> >` : une liste par ligne de la matrice, chaque liste contient les indices des colonnes non nulles → dynamique mais pas de contiguïté
- ▶ `vecteur<list<unsigned int> >` : si on connaît le nombre de lignes, comme précédemment mais plus efficace ?
- ▶ `vecteur<set<unsigned int> >` : comme précédemment mais on a le tri et l'unicité des indices des colonnes → le tri est-il nécessaire à chaque étape ?

└ Structures de données

└ Liste triée - set<.>

---

**Exple :** Une matrice creuse en 0/1 (cas d'école ici) est une liste d'éléments non nuls avec leur index. Comment l'interroger ? A-t-on besoin de la modifier ?

- ▶ `vector<list<unsigned int> >` puis `sort/uniq` puis `vector<vector<unsigned int> >` : on passe par une structure temporaire pour aller vers la structure la plus efficace → poids de la copie versus poids de l'utilisation de la structure finale

Quid des performances en machine : il faut tester/profiler.

Les conteneurs STL sont hyper-optimisés, les différences sont donc sensibles au delà d'une certaine taille de problèmes.

C++ coding standards, Herb Sutter, item 8 "Don't optimize prematurely"

## Structures de données

### Dictionnaire et table de hachage - `map<.,.>` et `unordered_map<.,.>`

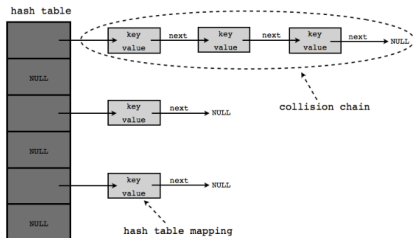
Associer des valeurs à des clés qui ne sont pas des entiers consécutifs  
(containeur associatif)

`map`

- ▶ `map<.,.> == set<pair<.,.> >`
- ▶ access avec opérateur `[ ]` mais attention  $O(\log n)$  !
- ▶ même comportement que le `set<.>` pour le reste

`unordered_map`

- ▶ access par fonction de hachage  
(calcul d'un index depuis une clé) en  $O(1)$  au mieux
- ▶ comme son nom l'indique, les éléments ne sont pas ordonnés mais on peut les parcourir
- ▶ très efficace en pratique



## └ Structures de données

## └ Dictionnaire et table de hachage - map&lt;.,.&gt; et unordered\_map&lt;.,.&gt;

```
1 #include <map>
2 #include <vector>
3 #include <string>
4 #include <stdlib.h>
5 using namespace std;
6 int tirage(int n){
7     return (int)((double)rand()/((double)RAND_MAX + 1) * n);
8 }
9 int main()
10 {
11     vector<string> v;
12     map<int, string> m;
13     int n = 1e7;
14
15     /***** init *****/
16     for (int i=0; i<n; i++)
17         v.push_back("PuyDeDome");
18     for (int i=0; i<n; i++)
19         m[i] = "PuyDeDome";
20
21     /***** acces *****/
22     // 4 secondes
23     for (int i=0; i<n; i++){
24         int j = tirage(n);
25         v[j] = "Cantal";
26     }
27
28     // 25 secondes
29     for (int i=0; i<n; i++){
30         int j = tirage(n);
31         m[j] = "Cantal";
32     }
33 }
```

└ Structures de données

└ Dictionnaire et table de hachage - `map<.,.>` et `unordered_map<.,.>`

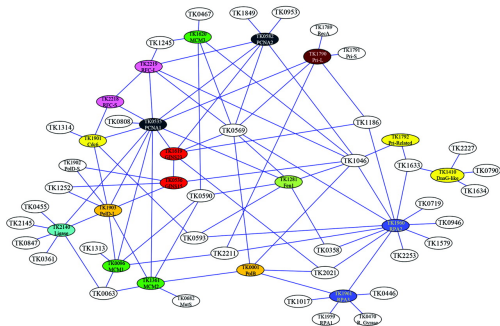
---

Exple : Une matrice creuse à valeurs réelles...

- ▶ `map<pair<unsigned int, unsigned int>, double>` : pratique mais pas de contiguité en mémoire, empreinte mémoire conséquente
- ▶ triplet de `vector` : les valeurs, le positionnement par ligne et les indices de colonnes → encapsulation
- ▶ s'appuyer sur une librairie



Exple :



Nécessité de retenir les noms des  $n$  sommets d'un graphe pour le pre/post-traitement, mais inutilité dans les étapes de calculs :

- ▶ conversion des noms en entiers consecutifs + stockage de arêtes en matrice creuse
- ▶ utilisation de `map<string, unsigned int>` pour la conversion ou `vector<string>` + `find` en  $O(\log(n))$
- ▶ attention à la consommation de mémoire additionnelle

Use vector by default !

**Lecture** : Effective STL, Scott Meyers, item 23 “Consider replacing associative containers with sorted vectors”

**Lecture** : C++ coding standards, Herb Sutter, item 76 “Use vector by default. Otherwise choose an appropriate container”

No container independent code !

**Lecture :** Effective STL, Scott Meyers, item 44 “Prefer member functions to algorithms with the same names”

“Be curious” = Lire les specs des structures proposées (cppreference, Wikipedia, . . . )

**Lecture :** C++ coding standards, Herb Sutter, item 8 “Don’t optimize prematurely”

**Lecture :** The C++ standard library, Nicolai Josuttis, page 10, “Remember to think big when you consider complexity”

Faire des benchmarks pour se convaincre.

“Be lazy” = Toujours utiliser les structures de données des bibliothèques au plus haut niveau possible.

└ Se documenter

└ Le bon reflexe

---

Ouvrir des livres et des revues

+ faire des benchmarks et du profiling (voir cours suivants)

+ participer à des réseaux métiers/collègues

+ éviter tout dogmatisme

+ **Lecture** : C++ coding standards, Herb Sutter, item 6 "Correctness, simplicity and clarity come first"