



INFO2050-1 - Programmation Avancée

Projet 2 : Structures de Données

Aliaksei Mazurchyk

Antoine Sadzot

26 novembre 2017

Université de Liège

1 Analyse Théorique

1.1 Variante d'union-find à base d'arbres choisie

La version de l'union-find en arbre implémentée est une version améliorée. Comme dans la version de base chaque élément pointe vers son parent. Lors d'une union de sous-ensembles, les racines des deux arbres sont fusionnés, un parent se retrouve donc avec deux enfants.

1.1.0.1 Rang La première amélioration consiste à associer un rang à chaque élément. Cela permet d'optimiser la fusion. L'union-find de base pouvait être fort déséquilibré, ce qui rendait l'appel `Find()` assez long. L'ajout d'un rang permet d'équilibrer l'arbre lors de l'union.

Au début tous les éléments ont un rang 0. Lors d'une fusion de deux arbres de même hauteur, le rang de la nouvelle racine est incrémenté de 1. Si les rangs sont différents, l'arbre ayant le rang le plus petit est ajouté à la racine de l'arbre le plus grand. Cela permet d'équilibrer la structure au fur et à mesure de la construction et ainsi de diminuer la hauteur maximale.

1.1.0.2 Compression de chemin Cette seconde amélioration consiste à utiliser la fonction `Find()` pour aplatir l'arbre. Cela se fait lors de l'appel récursif, la racine de l'arbre remonte dans la call-stack et est définie en tant que parent de chaque neud parcouru. De cette manière, tous les éléments parcouru par un `Find()`, pointeront directement vers la racine et nécessiteront donc un temps constant pour récupérer cette dernière.

1.2 Complexités en temps de `ufUnion` et `ufFind`

1.2.1 Implémentation à base de listes chaînées

ufFind : Dans tous les cas, la complexité est $\Theta(1)$

ufUnion : Dans le meilleur cas, le plus petit ensemble à lier à la suite du second ensemble ne contient qu'un élément. La complexité dans le meilleur cas est donc $\Theta(1)$. Dans le pire cas, l'ensemble à lier à la suite du second ensemble est de taille $N/2$, N étant le nombre d'éléments du labyrinthe. La complexité dans le pire cas est donc $\Theta(N/2)$.

1.2.2 Implémentation à base d'arbres

ufFind : Dans le meilleur des cas, la complexité est $\Theta(1)$. Dans le pire est est $\mathcal{O}(h)$, avec h comme auteur de l'arbre.

ufUnion : Dans le meilleur des cas, la complexité est $\Theta(1)$. Dans le pire est est $\mathcal{O}(n + \log(n))$.

Complexité combinée : Dans le meilleur des cas, la complexité est $\Theta(1)$. Dans le pire elle est $\Omega(\alpha(n))$, ce qui équivaut à un temps constant car $\alpha(n)$ (fonction d'Ackermann inversée) a une croissance extrêmement lente. En pratique : $\alpha(2^{65536}) = 5$ ce qui garanti une complexité $\Theta(5)$. Cette complexité vaut pour l'*Union* et le *Find*, car les deux améliorations sont complémentaire et l'un tire parti des améliorations de l'autre. Ceci fut démontré par Fredman et Saks en 1989.

1.2.2.1 Sources

- en.wikipedia.org/wiki/Disjoint-set_data_structure
- fr.slideshare.net/NikitaShpilevoy/algorithms-union-find
- fr.slideshare.net/WeiLi73/time-complexity-of-union-find-55858534
- github.com/emintham/Papers/blob/master/Fredman%2CSaks-%20The%20cell%20probe%20complexity%20of%20dynamic%20data%20structures.pdf

1.3 Implémentation de la structure du labyrinthe

Les données du labyrinthe sont stockées dans la structure Maze.

```
1 struct maze_t {
2     UnionFind *union_Find;
3     struct Walls *myWalls;
4     size_t size;
5     size_t number_inner_walls;
6 };
```

La structure "maze_t" contient la structure "UnionFind", qui sera décrite pour chaque implémentation par la suite. Elle contient également la structure "Walls" qui est un vecteur enregistrant tous les murs internes du labyrinthe. La variable "size" contient la taille du labyrinthe. La variable "number_inner_walls" contient le nombre de mur internes du labyrinthe. C'est également la taille du vecteur de type "Walls".

```
1 typedef struct Walls {
2     Coord Cell1;
3     Coord Cell2;
4     bool wall_between;
5 };
```

La structure "Walls" contient les coordonnées de deux cellules adjacentes, partageant un mur. La variable booléenne "wall" vaut "true" si le mur est présent. Elle vaut "false" si le mur est absent.

1.3.1 Implémentation par listes chaînées

La structure union_find_t est décrite de cette manière :

```
1 struct union_find_t {
2     struct Element* elements;
3     struct Sentinel* sentinels;
4     size_t numberComponents;
5 };
```

Elle contient un vecteur "elements" contenant tous les éléments du labyrinthe structurés en listes chaînées. Elle contient également un vecteur "sentinels" contenant toutes les sentinelles pointant sur les premiers et derniers éléments de chaque liste chaînée. La variable "numberComponents" contient le nombre d'ensembles distincts.

```
1 typedef struct Element {
2     size_t numero;
3     struct Sentinel* head;
4     struct Element* next;
5 }Element;
```

Le "numero" correspond à l'indice de l'élément dans le labyrinthe. Le pointeur "head" pointe vers la sentinelle de l'ensemble auquel appartient l'élément. Le pointeur "next" pointe vers l'élément suivant dans la liste chaînée.

```
1 typedef struct Sentinel {
2     struct Element* first;
3     struct Element* last;
4     size_t numberElements;
5 }Sentinel;
```

Le pointeur "first" pointe vers le premier élément de la liste chaînée liée à la sentinelle. Le pointeur "last" pointe vers le dernier élément de cette liste. La variable "numberElements" contient le nombre d'éléments présents dans cette liste.

1.3.2 Implémentation par arbres

La structure union_find_t est déclarée ainsi :

```

1 struct union_find_t {
2     size_t* items;
3     size_t* parents;
4     size_t* rank;
5     size_t n_items;
6     size_t n_trees;
7 };

```

Cette structure est composée de 3 tableaux :

- **items[n]** : Indices dans la grille du labyrinthe.
- **parents[n]** : Parent direct de l’item située à l’indice n .
- **rank[n]** : Rang de l’item de l’item située à l’indice n

Deux autres variables sont définies :

- **n_items** : Nombre total d’éléments.
- **n_trees** : Nombre d’ensembles.

Cette structure simplifie grandement le parcourt des ensembles disjoints, les pointeurs étant masqués derrière des tableaux.

1.4 Pseudocode des fonctions `mzCreate` et `MzIsValid`

mzIsValid(maze)

```

1 if UF_COMPONENTS_COUNT(maze->union_Find) > 1
2     return false;
3 else
4     return true;

```

Par soucis de simplicité, certaines séquences d’instructions évidentes ont été simplifiées en appels de fonctions. Comme par exemple l’initialisation des murs.

mzCreate(size)

```

1 innerWalls = NUMBER_INNER_WALLS(size)
2 myMaze->size = size
3 myMaze->union_find = UF_CREATE(size * size)
4 myMaze->number_inner_walls = innerWalls
5 INITIALIZE_WALLS(myMaze->myWalls, size * size)
6 parcours = SETUP_PARCOURS_ORDER(myMaze->number_inner_walls)
7 wallsToTest = 0
8 while !MZ_IS_VALID(myMaze) && wallsToTest < innerWalls
9     visits = parcours[wallsToTest]
10     indexCell1 = INDEX_FROM_COORD(myMaze->myWalls[visits].Cell1)
11     indexCell2 = INDEX_FROM_COORD(myMaze->myWalls[visits].Cell2)
12     status = UF_UNION(myMaze->union_Find, indexCell1, indexCell2)
13     if status == UF_MERGED
14         close = MZ_IS_WALL_CLOSED(myMaze, myMaze->myWalls[visits].Cell1,
15 myMaze->myWalls[visits].Cell2)
16         MZ_SET_WALL(myMaze, myMaze->myWalls[visits].Cell1,
17 myMaze->myWalls[visits].Cell2, close)
18     wallsToTest++
19 return myMaze

```

1.5 Analyse de la complexité en temps de `mzCreate` et `mzIsValid`

1.5.1 Implémentation à base de listes chaînées

mzIsValid : La fonction `ufComponentsCount` s’écrit comme telle :

ufComponentsCount(union_Find)

1 **return** *union_Find* - > *numberComponents*

Dans tous les cas, sa complexité est $\Theta(1)$.

mzCreate : Seules les parties de la fonction présentant une complexité égale ou supérieure à N^2 sont présentées. Les autres sont ignorées car négligeables. **Dans tous les cas.** La ligne 3 est $\Theta(N^2)$ car chaque cellule est visitée une à une. A la ligne 5, N^2 itérations sont nécessaires car toutes les cellules sont visitées. Donc la complexité est $\Theta(N^2)$.

A la ligne 6, chaque mur intérieur est visité, ce qui fait $2N^2 - 2N$ opérations, donc $\Theta(N^2)$. La boucle à la ligne 8 va varier en fonction du meilleur et du pire cas.

Dans le meilleur cas, chaque itération permet une connexion de cellules. Cela fait donc $N^2 - 1$ itérations nécessaires. Au sein de la boucle, toutes les instructions ont une complexité linéaire. C'est également le cas pour la ligne 11 dans le meilleur des cas. Car à chaque connexion d'ensembles, seul un élément est connecté à la fois. Donc, **dans le meilleur cas**, la complexité en temps est de l'ordre de $\Omega(N^2)$.

Dans le pire cas, tous les éléments du vecteur *myWalls* sont visités. Cela fait donc $2N^2 - 2N$ itérations. Il n'y aura cependant toujours que $N^2 - 1$ opérations de connexion de cellules. Donc seule la ligne 11 aura une complexité variable. L'analyse de la complexité du *ufUnion* va se faire pour les cas où N est une puissance de 2, par facilité. Le pire cas correspond à celui où les ensembles connectés sont toujours de taille égale. Au début, tous les ensembles seront de taille 1, puis de taille 2, 4, 8, ... jusqu'à ce qu'il ne reste plus qu'un ensemble. En prenant le problème à l'envers, on part d'un ensemble de N^2 éléments à N^2 ensemble d'un éléments :

$$T(1) = 0T(m^2) = \frac{N^2}{2} + T\left(\frac{m^2}{2}\right) \text{ pour } m = N \quad (1)$$

Par Plug-and-Chug, on obtient la solution suivante :

$$\log_2(N^2) * \frac{N^2}{2} \quad (2)$$

Ce qui est la complexité totale de la ligne 11. Le pire cas de la boucle qui démarre en ligne 8 influencera donc la complexité totale. **Dans le pire cas**, la complexité est de l'ordre de $\Omega(\log_2(N^2) * N^2)$ On observe que la complexité en temps du pire et du meilleur cas sont proches, à un multiple de N^2 près.

2 Analyse empirique

2.1 Tests de performance

FIGURE 1 – List

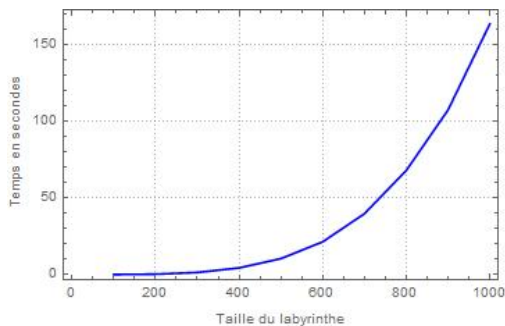


FIGURE 2 – Tree

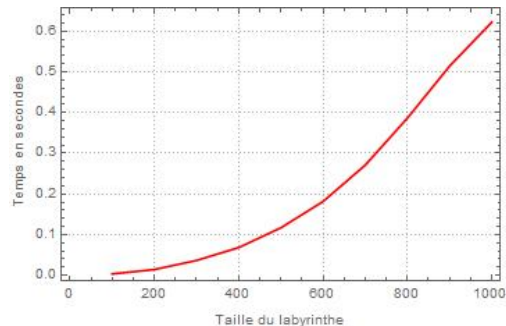


FIGURE 3 – List et Tree

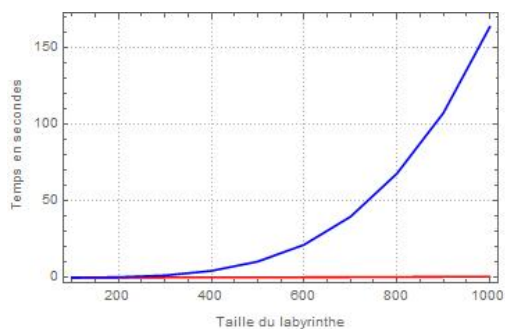
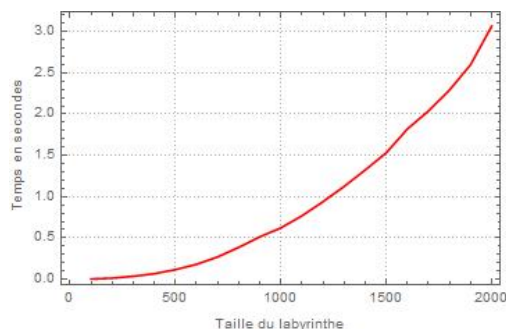


FIGURE 4 – Tree grande taille



2.2 Analyse des résultats

La figure 1 montre que le temps de calcul croît très rapidement avec le carré de la taille du labyrinthe. Ceci semble concorder avec les analyses théoriques combinées de *unionFindList* et de *mzCreate*.

Sur la figure 2 on peut remarquer que l'implémentation de *unionFindTree* croît très lentement, presque linéairement à partir de la taille 800. Des tests plus nombreux furent réalisés ce qui donna la figure 4. On y voit que la courbe croît très lentement. Le temps dédié aux opérations Union et Find devant être constant en pratique (démonstré en 1989 par Fredman et Saks), on peut supposer que la légère croissance vient des opérations de *mzCreate*.

La comparaison entre l'implémentation *List* et *Tree* illustrée par la figure 3 montre que la structure en arbre est bien plus performante. Pour un tableau de taille 1000, 10^6 éléments sont gérés par l'Union-Find. La génération du labyrinthe en utilisant la structure en liste nécessite 160 secondes, alors que celle en arbre ne prend que 0.6 secondes. Ceci est principalement dû aux améliorations implémentées à la structure en arbre (rang et compression de chemin).