

GESTION D'ENSEMBLES DISJOINTS

APPLICATION À L'ALGORITHME DE KRUSKAL

Ce problème est aussi connu sous le nom de **UNION - APPARTENANCE**, ou **UNION - FIND**, ou encore de **MERGE - FIND**, il s'agit en gros de gérer des ensembles disjoints à partir de n éléments indépendants

I - LE PROBLÈME

Étant donnés n éléments distincts notés $1, 2, \dots, n$, le problème consiste à gérer le regroupement de certains d'entre eux en sous-ensembles de plus en plus gros, que nous appellerons **CLASSES**

Initialement, il y a n classes : $\{1\}, \{2\}, \dots, \{n\}$

Chaque étape correspond au regroupement de deux éléments, x et y par exemple, ce que l'on notera $x \equiv y$, il s'agit alors de déterminer les classes C_x et C_y de x et y , puis d'en faire l'union $C_x \cup C_y$, si elles sont différentes. La suite des classes évolue donc au cours de regroupements, mais les sous-ensembles correspondants restent disjoints, on peut donc les assimiler à des classes d'équivalence

Par exemple, si $n = 7$, on a au départ $\{1\}, \{2\}, \{3\}, \{4\}, \{5\}, \{6\}$ et $\{7\}$. Les regroupements successifs de 1 et 2, puis de 5 et 6, puis de 3 et 4 et enfin de 1 et 4 donneront finalement les classes $\{1, 2, 3, 4\}, \{5, 6\}$ et $\{7\}$

Le problème est donc tout d'abord d'initialiser les classes avec les n éléments, nous appellerons cela **INITIALISATION**, il faut créer la classe $C_x = \{x\}$ pour chaque x . Puis à chaque étape de réaliser les opérations **APP(x)** et **APP(y)**, qu'on appellera **APPARTENANCE**, qui donnent les noms des classes auxquelles appartiennent x et y et **UNION** qui réalise l'union (ou la fusion) des classes C_x et C_y et place le résultat dans C_x ou dans C_y , il faut choisir, soit arbitrairement, soit suivant un certain critère, voir suite

Et le schéma général de l'algorithme, pour une suite de regroupements $x \equiv y$, est donc :

```
pour chaque  $x \equiv y$  faire  
     $C_x = \text{APP}(x)$   
     $C_y = \text{APP}(y)$   
     $\text{UNION}(C_x, C_y)$ 
```

Plusieurs implémentations sont possibles pour réaliser ces opérations, nous en étudions quelques-unes

II - UNE PREMIÈRE SOLUTION

Une première solution consiste à utiliser un unique tableau **Classe(1), Classe(2), ..., Classe(n)** où **Classe(x)** est un entier de $\{1, 2, \dots, n\}$ représentant la classe de x .

Initialement, on prendra **Classe(x) = x** pour tout élément, **INITIALISATION** est donc réalisée par :

```
pour  $x = 1$  à  $n$  faire  $\text{Classe}(x) = x$ 
```

Pour chaque $x \neq y$, **APPARTENANCE** consiste à consulter simplement $\text{Classe}(x)$ et $\text{Classe}(y)$, donc **UNION**, qui réalise $C_x = C_x \cup C_y$, peut s'écrire :

```
si  $\text{Classe}(x) \neq \text{Classe}(y)$  alors
  pour  $i = 1$  à  $n$  faire
    si  $\text{Classe}(i) = \text{Classe}(y)$  alors  $\text{Classe}(i) = \text{Classe}(x)$ 
```

En effet on passe en revue tous les éléments et tout élément de C_y devient un élément de C_x

L'**INITIALISATION** nécessite n étapes élémentaires, ainsi que chaque **UNION**
Une suite de m regroupements sera donc en $O(m.n)$

III - UNE SOLUTION PLUS EFFICACE

Dans la solution précédente, pour chaque union, on examine systématiquement tous les éléments, même ceux qui ne sont pas concernés. Pour éviter ces consultations inutiles, il suffit, pour chaque classe, de stocker l'ensemble de ses éléments, et uniquement ceux là

Au tableau $\text{Classe}(1), \dots, \text{Classe}(n)$, on ajoutera un tableau de listes indicé par le nom des classes : **Liste(1)**, **Liste(2)**, ..., **Liste(n)** où $\text{Liste}(i)$ est un pointeur sur une liste contenant les éléments de la classe i

La classe d'un élément sera l'élément qui est choisi comme représentant de la classe. Initialement, on a bien sûr $\text{Liste}(i) = \{i\}$, et à chaque étape on pourra simplement prendre l'élément situé en tête de liste

Les opérations pour **INITIALISATION** et **UNION** sont alors respectivement :

```
pour  $i = 1$  à  $n$  faire
   $\text{Classe}(i) = i$ 
   $\text{Liste}(i) = \{i\}$ 

si  $\text{Classe}(x) = C_x \neq \text{Classe}(y) = C_y$  alors
   $\text{Liste}(C_x) = \text{Liste}(C_x) \cup \text{Liste}(C_y)$ 
   $\text{Liste}(C_y) = \emptyset$ 
  pour chaque élément  $z$  de  $\text{Liste}(C_y)$  faire  $\text{Classe}(z) = \text{Classe}(x)$ 
```

APPARTENANCE, comme précédemment, consiste tout simplement à consulter $\text{Classe}(x)$ et $\text{Classe}(y)$

INITIALISATION se fait en $\Theta(n)$ et **UNION** en $\Theta(|C_y|)$, pour éviter de parcourir inutilement $\text{Liste}(C_x)$, on rajoutera un pointeur sur la fin de liste. Cela améliore les performances en moyenne par rapport à la solution précédente, mais une suite de m regroupements pourra se faire en $O(m.n)$ dans le plus mauvais cas

L'intérêt de cette solution est que l'on peut encore améliorer cette complexité en faisant une **UNION pondérée** ou **différenciée** qui consiste à rajouter C_y à C_x quand $|C_x| \geq |C_y|$, ou l'inverse sinon. Il faut alors rajouter dans le tableau Liste une case où l'on stockera la taille de la liste correspondante

Ainsi, **UNION** pourra se faire en $O(\min(|C_x|, |C_y|))$

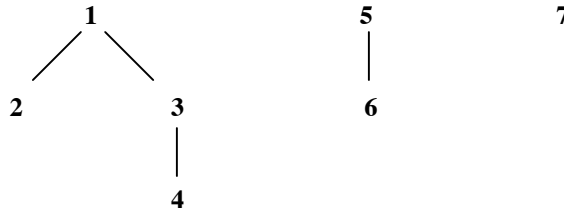
Tout élément changeant de classe se retrouve dans une classe de taille au moins double de celle où il se trouvait précédemment. Comme il y a initialement n classes, un élément ne peut changer de classe plus de $\log_2 n$ fois. Ainsi m regroupements pourront se faire en $O(m \log n)$

IV – FORÊT D'ARBORESCENCES

Cette méthode consiste à **représenter chaque classe par une arborescence dont la racine servira de représentant**, donc l'ensemble des classes sera une forêt

Dans cette solution, un seul tableau **Forêt(1), Forêt(2), ..., Forêt(n)** suffit à tout représenter : pour tout élément x , non racine, $\text{Forêt}(x) = y$ si y est le père de x dans une arborescence et $\text{Forêt}(x) = -k$ si x est racine, où k est le nombre d'élément d'éléments de l'arborescence

Avec l'exemple donné en introduction, on a :



Le tableau **Forêt** correspondant, indicé de 1 à 7, sera alors :

- 4	1	1	3	- 2	5	- 1
-----	---	---	---	-----	---	-----

INITIALISATION s'écrit alors :

pour $i = 1$ à n **faire** $\text{Forêt}(i) = - 1$

APP(x) consistera alors, partant de x , à « remonter » jusqu'à la racine :

$C_x = x$
tant que $\text{Forêt}(C_x) > 0$ **faire** $C_x = \text{Forêt}(C_x)$

UNION entre $C_x = \text{APP}(x)$ et $C_y = \text{APP}(y)$ peut se faire en rajoutant C_y comme fils de C_x , d'où :

$\text{Forêt}(C_x) = \text{Forêt}(C_x) + \text{Forêt}(C_y)$
 $\text{Forêt}(C_y) = C_x$

Mais là aussi on peut obtenir une composante réduite à une simple chaîne, une suite de m regroupements peut donc dégénérer en $O(m.n)$. Pour remédier à ce problème il suffit de tenir compte de la taille des arborescences et de rajouter la racine de la plus petite comme fils (fille ?) de la racine de la plus grande :

si $\text{Forêt}(C_x) \leq \text{Forêt}(C_y)$ **alors** $\text{Forêt}(C_x) = \text{Forêt}(C_x) + \text{Forêt}(C_y)$; $\text{Forêt}(C_y) = C_x$
sinon $\text{Forêt}(C_y) = \text{Forêt}(C_x) + \text{Forêt}(C_y)$; $\text{Forêt}(C_x) = C_y$

Dans l'exemple précédent, si le regroupement suivant est **4 = 6**, le tableau **Forêt** devient, car le sommet 5 est désormais fils de 1 :

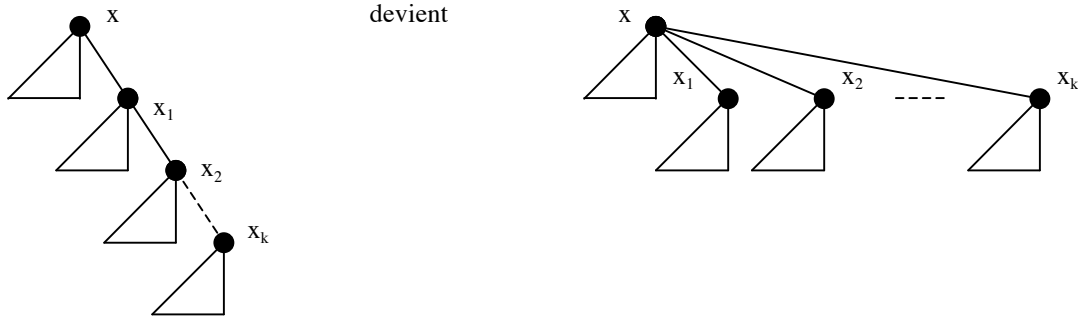
- 6	1	1	3	1	5	- 1
-----	---	---	---	---	---	-----

Alors une suite de m regroupements peut se faire en **$O(m \log n)$** , en effet, on peut montrer que la hauteur de toute arborescence est toujours au plus $\log_2 n$

V - COMPRESSION DE CHEMINS

On peut encore améliorer les performances en utilisant la « **Compression de Chemin** ». L'idée est la suivante : lors du parcours d'un élément x à la racine dans la fonction appartenance, chacun des sommets du chemin parcouru est rendu fils de la racine

Schématiquement :



Cette compression peut se faire lors de l'appel de la fonction **APP**(x), de la façon suivante :

$C_x = x$
tant que Forêt(C_x) > 0 **faire** Empiler(C_x) ; $C_x = \text{Forêt}(C_x)$
pour tout élément z de la pile **faire** Forêt(z) = C_x

On utilise une pile pour stocker les éléments du chemin allant de x à la racine

Cette compression « écrase » l'arborescence en hauteur qui aura pour effet d'accélérer la fonction d'appartenance. Les calculs de complexité sont difficiles, mais l'on a le résultat suivant :

Théorème

La compression de chemins avec l'union différenciée pour m regroupements est en $O(m \cdot \alpha(n))$

La fonction $\alpha(n)$ est définie sur les entiers par : $\alpha(n) = \min \{ x \text{ entier} \mid n \leq A(x, x) \}$ où A est la **Fonction d'Ackermann**, d'une croissance très rapide, qu'on peut définir pour tout couple d'entiers positifs par

$A(0, y) = 1$
 $A(1, 0) = 2$
 $A(x, 0) = x + 2$ pour $x \geq 2$
 $A(x, y) = A(A(x - 1, y), y - 1)$, pour $x, y \geq 2$

On peut montrer que :

- $A(x, 1) = 2x$, pour $x \geq 1$
- $A(x, 2) = 2^x$
- $A(x, 3) = \ll x \text{ élévations à la puissance de } 2 \gg$, pour $x \geq 1$
- $A(0, 4) = 1, A(1, 4) = 2, A(2, 4) = 2^2, A(3, 4) = 2^{16} = 65\,536, A(4, 4) = \ll 65\,536 \text{ élévations à la puissance de } 2 \gg, \dots$

Nous avons : $\alpha(1) = \alpha(2) = 1$, car $A(1, 1) = 2$; $\alpha(3) = \alpha(4) = 2$, car $A(2, 2) = 4$; $\alpha(5) = \dots = \alpha(16) = 3$, car $A(3, 3) = 16$, et $\alpha(n) \leq 4$ pour des valeurs « raisonnables » de n

α est donc une fonction de croissance extrêmement lente, qu'on peut considérer comme constante pour des tailles usuelles de données

VI – APPLICATION À L'ALGORITHME DE KRUSKAL

L'algorithme de **KRUSKAL** qui construit un **Arbre Recouvrant Minimal** d'un graphe simple non orienté valué $G = (X, E, v)$ procède en deux temps :

- Trier les arêtes par ordre de coût croissant
- Considérer les arêtes xy une à une dans cet ordre : toute arête ne formant pas de cycle avec les arêtes déjà retenues est gardée ; l'algorithme stoppe dès que l'on a $n - 1$ arêtes

La première étape peut se faire en temps en $O(m \cdot \log m) = O(m \cdot \log n)$, car $m \leq n^2$, tout le problème réside donc dans le test « d'**acyclicité** » de la deuxième étape : il faut tester si les deux sommets x et y de l'arête à tester sont dans le même sous arbre, autrement dit dans la même composante connexe de la forêt en construction. Cela peut coûter globalement $O(m \cdot n)$, soit $O(n^3)$, complexité moins bonne que celle de la première étape

Il est possible de faire bien mieux en remarquant que la deuxième étape est exactement un problème de type **UNION – FIND**, il suffit en effet de noter que l'on a à gérer des classes d'équivalence, chaque classe correspondant aux sommets des composantes connexes de la forêt. Ainsi les méthodes développées précédemment peuvent s'appliquer

Pour finir, remarquons qu'il faut aussi stocker l'arbre en construction, une simple liste d'arêtes peut suffire, ou, mieux, un tableau de listes d'adjacence, l'espace sera en $\Theta(n)$, dans les deux cas. Stocker l'arbre dans une matrice d'adjacence n'est pas judicieux, celle-ci nécessitant un temps et un espace en n^2 , pour une structure d'arbre n'ayant que $n - 1$ arêtes

En conclusion, en prenant en compte le stockage du graphe G en tableau de listes d'adjacences :

- Le Tri des arêtes peut se faire en $\Theta(m \cdot \log n)$ en temps (Tri par Fusion, par exemple) et en $\Theta(n + m)$ en espace
- La construction de l'Arbre peut se faire en $O(m \cdot \alpha(n))$ en temps en $\Theta(n + m)$ en espace

Si G est donné par sa matrice d'adjacence, la complexité en espace sera en $\Theta(n^2)$