

TP1 - Algorithme probabiliste pour la coupe minimum

OPTIMISATION ET RECHERCHE OPERATIONNELLE

M1 Info - semestre d'automne 2021-2022

UNIVERSITÉ CLAUDE BERNARD LYON 1

Christophe Crespelle

christophe.crespelle@inria.fr

Le but de ce TP est d'implémenter, en langage C, l'algorithme probabiliste pour la coupe minimum vu en cours. Plusieurs fichiers sont fournis : le fichier contenant le programme principal `main.c` et trois bibliothèques `graph.h`, `nodelist.h` et `utility.h`, avec les `.c` correspondant. Le fichier `main.c` contient une trame de l'architecture générale du programme, à compléter, avec des sections pré-définies en commentaire qui vous indiquent où placer certaines des parties de code que vous allez développer. La bibliothèque `graph.h` contient des fonctions de manipulation des graphes et `nodelist.h` contient des fonctions de manipulation de listes de nœuds. Vous coderez les fonctions nécessaires à l'algorithme dans une quatrième bibliothèque que vous créerez et que vous inclurez de la même manière que les autres au début du `main.c`.

Le programme se compile grâce à la commande `make` tapée dans un répertoire contenant le fichier `Makefile` fourni ainsi que tous les fichiers du programme. Il s'exécute en tapant `./mincut.out`. Tapez `./mincut.out -h` pour afficher l'aide.

1 Prise en main du programme

Question 1. Compilez le programme. Combien d'options a-t-il ? À quoi servent-elles ? Quelles sont les valeurs par défaut ?

Question 2. Quelles fonctions contient la bibliothèque `graph.h` ? À quoi sert la fonction `graph_from_file` ?

Dans quelle structure est stocké un graphe ? De quel type de graphe s'agit-il ?

Question 3. Quelles fonctions contient la bibliothèque `nodelist.h` ? Quelle structure stocke une liste d'adjacence ? De quel type de liste s'agit-il ?

Pour l'instant, le programme se contente d'ouvrir les fichiers d'entrée et de sortie spécifiés lors de l'appel au programme, respectivement en lecture seule et écriture seule, et de les fermer à la fin du programme.

Question 4. En utilisant les fonctions de la bibliothèque `graph.h`, complétez le `main` pour que le programme :

- charge un graphe G en mémoire à partir du fichier spécifié en entrée (ou de l'entrée standard),

- ecrive sur la sortie d'erreur standard (`stderr`) le nombre de sommets et le nombre d'aretes du graphe charge,
- libere la place prise par le graphe en memoire avant de quitter.

Vous fermerez le fichier d'entree immediatement apres avoir charge le graphe qu'il contient, a l'aide de la fonction `fclose`, deja utilisee a la fin du main.

Indication. Vous trouverez de nombreux exemples d'écriture dans un fichier ou sur les sorties standard dans les fichiers qui vous sont fournis. L'annexe B du sujet contient également un petit rappel sur l'utilisation de la fonction `fprintf`.

Plusieurs fichiers qui contiennent chacun un graphe, et dont le nom est de la forme `graphEL_*`, vous sont donnees dans le dossier `/data` accessible depuis la page web des TP (voir descriptif succinct de leur provenance dans la section 3). Le format qu'ils utilisent pour encoder le graphe G est le suivant :

- la premiere ligne du fichier contient le nombre n de sommets dans le graphe,
- toutes les autres lignes du fichier contiennent une arete du graphe sous la forme de deux entiers u v separees par un espace, ou u et v sont les identifiants des deux sommets (distincts) qui sont les extremités de l'arete,
- l'identifiant d'un sommet est un entier compris entre 0 et $n - 1$.

Question 5. Essayez votre programme sur quelques-uns des graphes fournis, en particulier sur `toygraph` et `toygraph2`. La sortie de votre programme est-elle correcte ?

2 Implementation de l'algorithme

Vous implementerez toutes les fonctions necessaires pour l'algorithme dans une bibliotheque a part, nommee par exemple `fonctions-algo.h`, que vous appellerez au debut de `main.c` par un `#include`.

2.1 Contracter une arete

Le coeur de l'algorithme vu en cours est la fonction de contraction d'une arete. Une des implementations les plus simples de cet algorithme consiste a simuler les contractions d'aretes sans les effectuer reellement sur le graphe, c'est a dire que le graphe reste inchangé durant tout l'algorithme. Pour cela on attribue a chaque sommet l'identifiant du groupe auquel il appartient, qui est toujours choisi comme l'identifiant d'un sommet du groupe. Initialement chaque sommet est dans un groupe different (dont l'identifiant est donc l'identifiant du sommet lui meme) et lors d'une contraction d'arete, on se contente de faire l'union des groupes des deux sommets qui sont les extremités de l'arete contractee. A la fin de l'algorithme, lorsqu'il ne reste que deux groupes, on obtient la valeur de la coupe correspondante en comptant le nombre d'aretes entre deux sommets qui n'appartiennent pas au meme groupe.

En terme de structure de donnee, on representera les groupes a l'aide de trois structures :

- un tableau de taille n qui a chaque sommet associe l'identifiant de son groupe et
- un tableau de taille n qui a chaque identifiant de groupe associe la liste des sommets qu'il contient et

— un tableau de taille n qui a chaque identifiant de groupe associe le nombre de sommets dans le groupe.

Question 6. Dans la section *DATA STRUCTURE* du `main`, declarez les trois tableaux representant les groupes, reservez l'espace memoire necessaire pour chacun d'eux et initialisez les comme ils doivent l'etre au debut de l'algorithme.

Indication. Avant de vous lancer, lisez (ou relisez, si besoin) l'annexe A sur la gestion de la memoire en C.

Pour des questions de complexite, lors d'une contraction, on changera seulement les identifiants des sommets du plus petit groupe. Le probleme d'union des groupes auquel on se retrouve confronte est un probleme tres classique en informatique connu sous le nom d'*union-find*.

Question 7. Vous avez du deja voir le probleme d'union-find au cours de vos etudes. Quelle est la complexite totale de l'implementation proposee ici pour l'ensemble des unions realisees au cours de l'algorithme de contraction ?

Question 8.

Implementez la fonction `contraction_simulee` de la bibliotheque `fonctions-algo.h`, qui actualise, lors d'une contraction d'arete, les trois tableaux representant les groupes. Pensez a inclure cette bibliotheque au debut du fichier `main.c`.

Question 9. Dans la section *COMPUTE MIN CUT* du `main`, testez votre fonction, et verifiez que son resultat est correct, pour la contraction d'une seule arete, puis pour quelques contractions successives. Faites cela avec l'un des deux graphes jouets fournis dans le dossier `/data` accessible depuis la page web des TP. Vous pourrez utiliser la fonction `write_graph` de la bibliotheque `graph.h` pour verifier le graphe contracte obtenu.

Indication. Lors des contractions successives, prenez garde a ne contracter que des aretes qui existent encore dans le multigraphe contracte courant, c'est a dire qui sont entre des sommets appartenant a des groupes differents.

2.2 Tirer une arete au hasard

A chaque etape de l'algorithme, on doit tirer uniformement aleatoirement une arete parmi les aretes du multigraphe \tilde{G} courant (obtenu par contractions successives du graphe de depart), que nous avons choisi de ne pas explicitement represente. Une facon simple et efficace de faire ces tirages aleatoires est de choisir d'emblee un ordre aleatoire π sur toutes les m aretes du graphe G de depart. Au cours de l'algorithme on parcourt cet ordre en considerant les aretes une par une : si l'arete consideree est une arete du multigraphe \tilde{G} courant, c'est celle-ci que l'on contracte a cette etape de l'algorithme, sinon, si l'arete consideree a disparu du multigraphe dans les operations de contraction effectuees precedemment dans l'algorithme, on la rejette et on passe a l'arete suivante dans l'ordre π .

Question 10. Cette facon de proceder garantie-t-elle que l'arete tiree a une etape de l'algorithme est choisie uniformement aleatoirement parmi les aretes restantes dans le multigraphe \tilde{G} ?

Question 11. Comment déterminer si une arête uv de G appartient au multigraphe contracté \tilde{G} courant, à l'aide des groupes auxquels appartiennent les sommets u et v ?

Pour stocker un ordre sur les arêtes, deux choix naturels sont de :

- soit utiliser un tableau de structures de type `arete` (un couple de sommets)
- soit utiliser deux tableaux de sommets, un contenant la première extrémité de l'arête, l'autre la deuxième extrémité.

Question 12. En utilisant la fonction `rand()` de tirage aléatoire du langage C, écrivez dans la bibliothèque `fonctions-algo.h`, une fonction qui génère un ordre aléatoire des arêtes du graphe G de départ.

Indication. Vous pourrez procéder comme suit. Commencez par stocker les arêtes du graphe dans le (ou les) tableau(s) résultat dans un ordre arbitraire, par exemple celui dans lequel vous rencontrez les arêtes en parcourant les listes d'adjacences du graphe G . Faites attention à ne stocker chaque arête qu'une seule fois (chacune est présente deux fois dans les listes d'adjacence). Ensuite, mélangez aléatoirement le tableau en choisissant uniformément aléatoirement à chaque étape i , pour i allant de 0 à $m - 1$, l'arête que vous placez en position i du tableau, parmi les arêtes qui sont dans les positions i à $m - 1$ (c'est à dire celles que vous n'avez pas encore tirées). Pensez à placer l'arête qui était en position i précédemment à la place de celle que vous avez tirée entre les positions i et $m - 1$.

2.3 Iteration de l'algorithme de contraction

Question 13. Écrivez la boucle qui contracte une par une les arêtes du graphe. Combien de fois cette boucle doit-elle s'exécuter ?

Question 14.

Dans la bibliothèque `fonctions-algo.h`, écrivez une fonction `valeur_coupe` qui calcule la valeur d'une coupe.

Pour avoir une bonne probabilité de trouver la coupe minimum, il faut appliquer l'algorithme de contraction aléatoire plusieurs fois, avec des tirages aléatoires d'arêtes différents.

Question 15. Dans le `main`, écrivez la boucle qui itère cet algorithme autant de fois que voulu, avec un ordre aléatoire sur les arêtes différents à chaque fois.

Question 16. Ajoutez le code permettant de retenir la plus petite coupe trouvée lors des itérations de l'algorithme probabiliste de contraction.

En fin de programme, on veut écrire dans un fichier résultat la coupe minimum trouvée. On adoptera le format suivant : la première ligne du fichier contiendra la valeur de cette coupe, la deuxième ligne la liste des identifiants des sommets qui sont dans la première partie de la coupe, séparés par un espace, et la troisième ligne la liste des identifiants des sommets qui sont dans la deuxième partie de la coupe.

Question 17. Écrivez à la fin du `main`, le code nécessaire pour écrire les résultats dans le fichier de sortie passé en paramètre du programme par l'utilisateur.

3 Application de l'algorithme sur des jeux de donnees

Plusieurs jeux de donnees synthetiques ou provenant de contextes reels sont donnees dans le dossier `/data` accessible depuis la page web des TP. L'annexe C en contient une breve description.

Question 18. Pour chacun des graphes fournis, verifiez le nombre de sommets ecrit dans le fichier et comptez le nombre d'aretes presentes dans le fichier.

Indication. Vous vous aiderez de la commande `head -1 fic` du systeme d'exploitation, qui affiche la premiere ligne du fichier `fic` et de la commande `wc -l fic` qui affiche le nombre de lignes dans le fichier `fic`.

Question 19. Appliquez l'algorithme, avec une unique iteration, sur chacun des graphes fournis. Quelle est la valeur de la coupe trouvee? Dans quel cas peut-on etre sur quil s'agit de la coupe minimum?

Question 20. D'apres le cours, comment choisir le nombre d'iterations N pour obtenir la coupe minimum avec une probabilite d'au moins 0,999? Pour quels graphes peut-on faire ce choix tout en gardant un temps d'execution raisonnable (c'est a dire avoir la reponse avant ce soir minuit)?

Indication. Utilisez la commande `time COMMANDE` du systeme d'exploitation qui donne le temps pris par l'execution de la commande `COMMANDE`. Faites cela pour chaque graphe en fixant le nombre d'iterations de l'algorithme a 10.

Question 21. En faisant un tres grand nombre d'iterations, c'est a dire en choisissant N comme indique a la question precedente, estimez la probabilite de tomber sur la coupe minimum en une seule iteration, pour chacun des graphes fournis. Vous considererez qu'au cours de ces N iterations, vous avez en effet decouvert la coupe minimum (verifiez en comparant au resultat obtenu par vos camarades).

Question 22. Pour les graphes qui sont trop volumineux pour choisir N comme demande a la question precedente, vous prendrez N le plus grand possible et en regardant les resultats obtenus vous discuterez de la validite de l'hypothese selon laquelle vous avez obtenu la coupe minimum.

A Rappel sur la gestion de la memoire en C

En C, on doit gerer soi-meme la memoire. On alloue de l'espace memoire sur le tas grace a la fonction `malloc` qui s'utilise ainsi :

```
TYPE* p;  
p = (TYPE*) malloc(sizeof(TYPE));
```

La premiere ligne declare un pointeur `p` sur un type `TYPE`, la deuxieme alloue sur le tas l'espace necessaire pour stocker une variable de type `TYPE`. L'argument de `malloc` est la taille en octet de la zone memoire allouee, le forçage de type `(TYPE*)` devant `malloc` est necessaire pour la concordance des types avec la variable `p` de type pointeur sur `TYPE`. Si on veut reserver un tableau de type `TYPE` et de taille `TAILLE`, on ecrit :

```
p = (TYPE*)malloc(TAILLE*sizeof(TYPE));
```

On accede alors a la case d'indice `i` du tableau `p`, pour `i` entre 0 et `TAILLE-1`, par l'expression `p[i]`.

Pour eviter les problemes lors du debogage, on protegera toutes les allocations memoires en faisant quitter le programme avec une message d'erreur a chaque fois qu'une allocation est ratee. Pour cela on utilisera la fonction `report_error` fournie dans la bibliotheque `utility.h`, de la maniere suivante :

```
if ( (p = (TYPE*) malloc(TAILLE*sizeof(TYPE))) == NULL )  
    report_error("malloc() error");
```

Vous avez de multiples exemples d'utilisation de `malloc` dans les fichiers qui vous sont fournis. Pour desallouer la zone memoire pointee par le pointeur `p` on utilise la fonction `free` :

```
free(p);  
p=NULL;
```

Il n'y a pas de garantie sur ce que vaut `p` apres la desallocation. C'est pour cela que par securite, il est conseille de mettre `p` a `NULL` apres desallocation pour eviter d'avoir des pointeurs vers des zones desallouees. Il n'y a pas non plus de garantie sur ce que contient la zone desallouee apres desallocation : elle peut etre reinitialisee ou lisee intacte (le plus courant) et dans tous les cas elle peut bien sur etre reservee plus tard dans le programme a un autre usage.

Les fuites memoires sont un vrai probleme, elles vous empecheront de traiter de grandes instances en entree de votre programme (ce qui est votre but). Il faut donc leur faire la chasse. Pour les eviter, vous desallouerez systematiquement toute memoire reservee, dans le programme principal (`main`) ou dans un sous-programme (fonction ou procedure), des lors que vous savez qu'elle ne servira plus (et dans tous les cas, au plus tard a la fin du `main`). Vous pouvez verifier que votre programme n'a pas de fuite memoire en scrutant la place qu'il utilise en memoire lors de son execution a l'aide de la fonction `top` du systeme d'exploitaion. Par exemple,

```
top -d 1 -u username -o %MEM
```

vous affiche les statistiques des programmes que vous (avec le nom d'utilisateur `username`) avez lance triees par ordre decroissant d'utilisation de la memoire (`-o %MEM`) et actualisees toutes les secondes (`-d 1`). Plus d'info avec `man top`.

B Rappel sur la fonction `fprintf`

La fonction `fprintf` permet d'écrire des données selon différents formats sur un flux de sortie. Un exemple d'utilisation simple est

```
fprintf(FLUX,"Texte quelconque");
```

qui écrit la chaîne de caractères "Texte quelconque" sur le flux de sortie `FLUX`, qui peut être un fichier (type `FILE*`) ouvert en écriture ou un flux de sortie standard du système d'exploitation, comme `stdout` (sortie standard) ou `stderr` (sortie d'erreur standard). Le flux est le premier paramètre de `fprintf` et la chaîne à écrire le deuxième paramètre. Il peut y avoir plus de paramètres dans l'appel à `fprintf`, placés après le flux et la chaîne, qui servent à incorporer d'autres données dans la chaîne de caractères en spécifiant le format auquel elles doivent être écrites sur le flux à l'aide du caractère `%`. Un exemple de syntaxe est le suivant :

```
fprintf(stderr,"Il y a %d sommets dans le graphe et %d arêtes",n,15);
```

Si `n` est une variable de type `int` qui vaut 8, cela écrit sur la sortie d'erreur standard la chaîne de caractères

```
Il y a 8 sommets dans le graphe et 15 arêtes
```

Le format d'écriture des entiers `n` et `15` est spécifié par le `%d` qui est le format d'écriture décimale des `int`. Les paramètres de `fprintf` qui suivent la chaîne de caractères, ici `n` et `15`, sont mis en correspondance avec les indications de format contenues dans la chaîne en utilisant l'ordre dans lequel apparaissent les paramètres supplémentaires passés à la fonction `fprintf` et l'ordre dans lequel apparaissent les indications de format dans la chaîne de caractères passée en paramètre à `fprintf`. Le nombre de paramètres supplémentaires doit donc être identique au nombre d'indications de format contenues dans la chaîne et il doit y avoir concordance entre les indications de format et les types des paramètres supplémentaires. Par exemple, `%u` est le format d'écriture décimale des `unsigned int`, `%lu` celui des `long unsigned int` et `%x` est le format d'écriture hexadécimale des `unsigned int`. On aurait par exemple pu appeler `fprintf` ainsi

```
fprintf(stderr,"Il y a %d sommets dans le graphe et %x arêtes",n,15);
```

ce qui écrit la chaîne de caractères

```
Il y a 8 sommets dans le graphe et f arêtes
```

Les retours à la ligne dans la chaîne de caractères sont codés par `\n`. Exemple :

```
fprintf(stderr,"Il y a %d sommets dans le graphe\net %d arêtes",8,15);
```

produit l'écriture

```
Il y a 8 sommets dans le graphe
et 15 arêtes
```

Pour optimiser les performances du système, l'écriture ne se fait pas directement sur le flux spécifié en paramètre de `fprintf` mais passe au préalable par un cache dont le contenu n'est effectivement écrit sur le flux que lorsque le cache est suffisamment rempli. On peut forcer l'écriture sur le flux et le vidage du cache à l'aide de la fonction `fflush` :

```
fflush(FLUX);
```

Afin de faciliter la phase de débogage, il est conseillé de procéder au vidage forcé du cache régulièrement, après chaque étape importante d'écriture. En effet, lorsque le programme quitte subrepticement, à cause d'un bug par exemple, le cache des flux ouverts en écriture n'est pas vide. Ainsi, tout ce qui avait été écrit sur le flux mais qui était encore dans le cache est perdu et n'apparaît pas sur le flux. Si l'on se fie à l'affichage sur le flux, cela trompe sur le moment de l'exécution auquel le bug s'est produit.

C Jeux de données

Dans le dossier `/data` accessible depuis la page web des TP, vous trouverez plusieurs graphes au format décrit au début du sujet. Certains de ces graphes sont synthétiques, c'est à dire qu'ils ont été générés soit à la main soit par des méthodes de génération automatiques, et d'autres proviennent d'opérations de mesure réalisées dans des contextes concrets : Internet, publications scientifiques, biologie, environnement, réseaux pair-à-pair, réseaux routiers et réseaux sociaux en ligne. Une brève description de ces graphes est donnée ci-dessous.

- `graphEL_toygraph`, 7 sommets, graphe jouet construit à la main,
- `graphEL_toygraph2`, 8 sommets, graphe jouet construit à la main,
- `graphEL_rand_100_8`, 100 sommets, graphe aléatoire (Erdos-Rényi G_{nm}) de degré moyen 8,
- `graphEL_rand_500_16`, 500 sommets, graphe aléatoire (Erdos-Rényi G_{nm}) de degré moyen 16,
- `graphEL_rand_1000_16`, 1000 sommets, graphe aléatoire (Erdos-Rényi G_{nm}) de degré moyen 16,
- `graphEL_as2000`, 6474 sommets, graphe de connections entre les systèmes autonomes d'Internet en 2000,
- `graphEL_ca-GrQc`, 4158 sommets, graphe de coauteurs dans le domaine de la relativité générale et de la cosmologie quantique,
- `graphEL_figeys`, 2217 sommets, graphe d'interactions chimiques entre protéines,
- `graphEL_foodweb`, 183 sommets, graphe de prédation entre espèces (chaîne alimentaire),
- `graphEL_p2p-Gnutella`, 62561 sommets, graphe de connections entre pairs dans le réseau P2P Gnutella,
- `graphEL_roadNet-TX`, 1351137 sommets, graphe du réseau routier du Texas,
- `graphEL_youtube`, 1134890 sommets, une partie du réseau social de youtube.