

TP1 - Algorithme probabiliste pour la coupe minimum

OPTIMISATION ET RECHERCHE OPERATIONNELLE

M1 Info - semestre d'automne 2020-2021

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. Deux fichiers sont fournis : `main.c` et `auxiliaire.c`, le second étant incorporé par un `#include` au début de `main.c`. Le fichier `main.c` contient une trame de l'architecture générale du programme (avec de sections pré-définies, en commentaire, qui vous indiquent où placer certaines des parties de code que vous allez développer) et le fichier `auxiliaire.c` contient des fonctions que vous utiliserez pour manipuler des graphes et leurs listes d'adjacences et tirer des nombres aléatoires. Vous coderez les fonctions nécessaires à l'algorithme dans un troisième fichier que vous créerez et que vous inclurez de la même manière au début de `main.c`.

Le programme se compile grâce à la commande `gcc -Wextra -Wall -O3 main.c -o mincut` et s'exécute en tapant `./mincut`. Tapez `./mincut -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. À quoi sert la fonction `graph_from_file` du fichier `auxiliaire.c` ? Quelle autre fonction de manipulation de graphe ce fichier contient-il ?

Question 3. Comment est représenté un graphe ? De quel type de graphe s'agit-il ?

Question 4. Quelle structure stocke une liste d'adjacence ? De quel type de liste s'agit-il ?

Question 5. Dans le fichier `auxiliaire.c`, écrivez une fonction `free_nod1` qui libère entièrement l'espace mémoire occupé par une liste de nœuds, donnée par un pointeur vers la liste. Attention, il faut libérer l'espace occupé par chacune des cellules de la liste !

De même, écrivez une fonction `free_graph` qui libère entièrement l'espace mémoire occupé par un graphe, la encore donnée par un pointeur sur le graphe.

Indication. Avant de vous lancer, lisez l'annexe A sur la gestion de la mémoire en C.

Pour l'instant, le programme se contente d'ouvrir les fichiers d'entree et de sortie specifiques lors de l'appel au programme, respectivement en lecture seule et ecriture seule, et de les fermer a la fin du programme.

Question 6. Completez le `main` pour que le programme :

- charge un graphe G en memoire a partir du fichier specifie en entree (ou de l'entree standard),
- ecrive sur la sortie d'erreur standard (`stderr`) le nombre de sommets et le nombre d'aretes du graphe charge,
- ecrive le graphe charge dans le fichier specifie en sortie (ou sur la sortie standard) et
- 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'ecriture dans un fichier ou sur les sorties standard dans les fichiers `main.c` et `auxiliaire.c` qui vous sont fournis. L'annexe B du sujet contient egalement 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 donnes 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 $u v$ separes 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 7. Essayez votre programme sur quelques-uns des graphes fournis, en particulier sur `toygraph` et `toygraph2`. Y a-t-il des differences entre le fichier d'entree et le fichier de sortie ecrit par la fonction `write_graph`? D'ou viennent ces differences?

Indication. Si besoin (pour les gros fichiers), vous pouvez comparer le contenu des deux fichiers de maniere automatique grace a la commande `diff fic1 fic2` du systeme d'exploitation.

Question 8. Avant de continuer, dans le `main`, gardez la partie du code qui charge le graphe en memoire et retirez celle qui l'ecrit dans le fichier de sortie (dans la suite, on va s'occuper d'y ecrire le resultat de l'algorithme).

2 Implementation de l'algorithme

Vous implementerez toutes les fonctions necessaires pour l'algorithme dans un fichier a part, nomme par exemple `fonctions-algo.c`, 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 inchange 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 9. 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 10. 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 11.

Dans le fichier `fonctions-algo.c`, implementez la fonction `contraction_simulee` qui actualise, lors d'une contraction d'arete, les trois tableaux representant les groupes. Incorporez le contenu du fichier `fonctions-algo.c` a l'aide d'un `#include "fonctions-algo.c"` au debut du fichier `main.c`.

Question 12. 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.

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, que nous avons choisi de ne pas explicitement représenter. 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 considerée 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 considerée 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 13. 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 14. Comment determiner si une arete uv de G appartient au multigraphe contracte \tilde{G} courant, a l'aide des groupes auxquels appartiennent les sommets u et v ?

Pour stocker un ordre sur les aretes, 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 premiere extremite de l'arete, l'autre la deuxieme extremite.

Question 15. En utilisant une des trois fonctions de tirage aleatoire contenu dans le fichier `auxiliaire.c`, ecrivez dans le fichier `fonctions-algo.c`, une fonction qui genere un ordre aleatoire des aretes du graphe G de depart.

Indication. Vous pourrez proceder comme suit. Commencez par stocker les aretes du graphe dans le (ou les) tableau resultat dans un ordre arbitraire, par exemple celui dans lequel vous rencontrez les aretes en parcourant les listes d'adjacences du graphe G . Faites attention a ne stocker chaque arete qu'une seule fois (elles sont presentes deux fois dans les listes d'adjacence). Ensuite, melangez aleatoirement le tableau en choisissant uniformement aleatoirement a chaque etape i , pour i allant de 0 a $m - 1$, l'arete que vous placez en position i du tableau, parmi les aretes qui sont dans les positions i a $m - 1$ (c'est a dire celles que vous n'avez pas encore tirees). Pensez a placer l'arete qui etait en position i precedemment a la place de celle que vous avez tiree entre les positions i et $m - 1$.

Indication. Afin de choisir laquelle des trois fonctions de tirage aleatoire fournies vous utiliserez, faite afficher les valeurs de `RAND_MAX`, `UINT_MAX` et `ULONG_MAX` **sur la machine sur laquelle vous travaillez**. Vous vous limiterez a traiter des graphes ayant au plus un milliard d'aretes.

2.3 Iteration de l'algorithme de contraction

Question 16. Ecrivez la boucle qui contracte une par une les aretes du graphe. Combien de fois cette boucle doit elle s'executer ?

Question 17. Dans le fichier `fonctions-algo.c`, ecrivez une fonction `valeur_coupe` qui calcule la valeur d'une coupe, donnee par les tableaux representant la 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 différentes.

Question 18. Écrivez la boucle qui itère cet algorithme autant de fois que voulu, avec un ordre aléatoire sur les arêtes différentes à chaque fois.

Question 19. 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 20. Écrivez à la fin du main, le code nécessaire pour écrire les résultats dans le fichier de sortie passe en paramètre du programme par l'utilisateur.

3 Application de l'algorithme sur des jeux de données

Plusieurs jeux de données synthétiques ou provenant de contextes réels sont donnés dans le dossier /data accessible depuis la page web des TP. L'annexe C en contient une brève description.

Question 21. Pour chacun des graphes fournis, vérifiez le nombre de sommets écrit dans le fichier et comptez le nombre d'arêtes présentes dans le fichier.

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

Question 22. Appliquez l'algorithme, avec une unique itération, sur chacun des graphes fournis. Quelle est la valeur de la coupe trouvée ? Dans quel cas peut-on être sûr qu'il s'agit de la coupe minimum ?

Question 23. D'après le cours, comment choisir le nombre d'itérations N pour obtenir la coupe minimum avec une probabilité d'au moins 0,999 ? Pour quels graphes peut-on faire ce choix tout en gardant un temps d'exécution raisonnable (c'est à dire avoir la réponse avant ce soir minuit) ?

Indication. Utilisez la commande `time COMMANDE` du système d'exploitation qui donne le temps pris par l'exécution de la commande `COMMANDE`. Faites cela pour chaque graphe en fixant le nombre d'itérations de l'algorithme à 10.

Question 24. En faisant un très grand nombre d'itérations, c'est à dire en choisissant N comme indiqué à la question précédente, estimez la probabilité de tomber sur la coupe minimum en une seule itération, pour chacun des graphes fournis. Vous considérez qu'au cours de ces N itérations, vous avez en effet découvert la coupe minimum (vérifiez en comparant au résultat obtenu par vos camarades).

Question 25. Pour les graphes qui sont trop volumineux pour choisir N comme demandé à la question précédente, vous prendrez N le plus grand possible et en regardant les résultats obtenus vous discuterez de la validité de l'hypothèse 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 un message d'erreur a chaque fois qu'une allocation est ratee. Pour cela on utilisera la fonction `report_error` fournie dans le fichier `auxiliaire.c`, 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 le fichier `auxiliaire.c`. 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 laissee 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 lancee 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

Vous trouverez plusieurs graphes au format décrit au début du sujet dans le dossier `/data` accessible depuis la page web des TP. 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 connection 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.