

INFO-F-413 : Implémentation de l'Algorithme de Karger

Thomas Chapeaux

September 30, 2012

1 Présentation

L'algorithme de Karger¹ est un algorithme probabiliste permettant d'estimer une coupe minimale d'un graphe, basé sur la contraction d'arêtes choisies aléatoirement.

Au cours, il a été vu que cet algorithme a une probabilité de succès (c'est-à-dire de trouver effectivement une coupe minimale) strictement supérieure à $\frac{2}{n(n-1)}$. Cette borne inférieure peut être améliorée jusqu'à $1 - \frac{1}{n^2}$ en exécutant $n(n-1) \log n$ itérations de l'algorithme et en gardant le meilleur résultat (la plus petite valeur trouvée).

Ce document présente une implémentation de l'algorithme en Python et l'utilise pour discuter les valeurs trouvées au cours.

2 Implémentation

2.1 Languages et bibliothèque

L'algorithme a été implémenté en Python 2.7 à l'aide de la bibliothèque `igraph`². Cette bibliothèque permet de générer des graphes aléatoirement via une méthode géométrique³ et propose une méthode pour trouver la coupé minimale d'un graphe, ce qui permettra de juger de la qualité des valeurs trouvées par notre algorithme.

Python a surtout été choisi par préférence personnelle, mais également pour son efficacité lors de projets de taille raisonnable comme celui-ci.

¹Réf. [1]

²<http://igraph.sourceforge.net>

³Réf. [2] et Annexe 1

2.2 Implémentation de l'algorithme

2.2.1 Une itération

Lors de la contraction d'une arête A, chaque autre arête connectée à la destination de A est redirigée vers la source de A⁴, et toutes les arêtes ayant la même source et la même destination que A (donc A également) sont supprimées.

```
1 # g is an instance of Graph()
2 while (g.vcount() > 2):
3     edge_list = g.get_edgelist()
4     rand_edge_id = random.randint(0, g.ecount()-1)
5     chosen_edge = edge_list[rand_edge_id]
6     chosen_edge_source = chosen_edge[0]
7     chosen_edge_target = chosen_edge[1]
8     for e in edge_list:
9         if (e[0] == chosen_edge_target or e[1] == chosen_edge_target):
10             if (e[0] == chosen_edge_target and e[1] != chosen_edge_source):
11                 :
12                 g.add_edge(chosen_edge_source, e[1])
13             elif (e[1] == chosen_edge_target and e[0] !=
14                   chosen_edge_source):
15                 g.add_edge(e[0], chosen_edge_source)
16             edge_id = graph.get_eid(e[0], e[1])
17             g.delete_edges(edge_id)
18
19     assert (g.degree(chosen_edge_target) == 0)
20     g.delete_vertices(chosen_edge_target)
```

2.2.2 Algorithme complet

On a choisi de faire par défaut le nombre d'itérations trouvé en cours pour avoir une probabilité de succès de $1 - \frac{1}{n^2}$ mais celui-ci peut être contrôlé en modifiant le paramètre *CONFIDENCE_FACTOR* pour faire moins d'itérations (facteur de confiance > 1) ou plus (< 1).

```
1 number_of_vertices = graph.vcount()
2 best_mincut_value = graph.ecount() + 1 # borne superieure
3 required_nbr_iter = int(number_of_vertices*(number_of_vertices-1)*math.
4   log(number_of_vertices))
5 nbr_iter = required_nbr_iter/CONFIDENCE_FACTOR
6 for i in range(nbr_iter):
7     g_temp = g.copy()
8     result = random_mincut_one_instance(g_temp)
9     if (result < best_mincut_value):
10         best_mincut_value = result
```

2.3 Outil d'analyse

2.3.1 Collecte de données

Le code va tester un certain nombre de graphes différents (générés aléatoirement). Lors de l'application de l'algorithme complet sur un graphe, certaines données sont collectées pour pouvoir en tirer des statistiques (voir plus loin).

⁴Le graphe étant non dirigé, on parle de source et de destination seulement du point de vue de l'implémentation dans *igraph*.

```

1 class TestResult:
2     def __init__(self, graph):
3         #self.graph = graph # not used
4         self.number_of_vertices = graph.vcount()
5         self.number_of_edges = graph.ecount()
6         self.number_of_iterations = None
7         self.improving_iterations = []
8         self.improved_value = []
9         self.found_value = None
10        self.correct_value = None
11
12    def consistencyCheck(self):
13        assert(len(self.improving_iterations)==len(self.improved_value))
14        assert(self.found_value != None)
15        assert(self.correct_value != None)

```

Une instance de cette classe par graphe testé est générée par le code pendant l'exécution de l'algorithme. On remarquera qu'on enregistre les itérations qui ont augmenté le résultat (*improving_iterations*) ainsi que la valeur correspondante (*improved_values*)

2.3.2 Analyse des données

Une fois tous les tests exécutés, le code tire les statistiques suivantes de tous les *TestResult* générés :

- Nombre de tests échoués (où une valeur trop grande a été trouvée pour la coupe minimale).
- Écart moyen avec la véritable valeur de coupe minimale.
- Proportions d'itérations "inutiles" (lorsque la valeur exacte avait déjà été trouvée)
- Itérations "utiles" les plus lointaines

Ces résultats sont ensuite affichés de la manière suivante :

```

1 confidence factor: 10
2 failures : 34 = 3.4% of tests
3 ...with an average of 1.32 absolute error
4 (Failed tests were removed from the following stats)
5 Three latest useful iterations (relative to total #iter) : 0.91, 0.9,
   0.8
6 Proportion of useful work : 0.15 (0.02 without confidence factor)

```

3 Discussion

Nous avons démontré (au cours) qu'en effectuant $n(n-1)\log n$ itérations, on était certain d'avoir un taux d'échec inférieur à $\frac{1}{n^2}$. Voyons comment ce résultat se traduit dans les tests.

Après une série de 1000 tests avec ce nombre d'itérations sur des graphes générés aléatoirement ayant au moins 5 sommets, nous devrions observer au plus $\frac{1000}{n^2} = 40$ erreurs. Le résultat est le suivant :

```

1 confidence factor: 1
2 failures : 0 = 0.0% of tests
3 Three latest useful iterations (relative to total #iter) : 0.4, 0.22,
  0.19
4 Proportion of useful work : 0.02 (0.02 without confidence factor)

```

On peut donc observer que la borne supérieure que nous avons trouvé est non seulement correcte, mais aussi fort exagérée⁵. Sur les 1000 tests, aucun n’a échoué et l’itération trouvant la bonne valeur la plus lointaine n’était qu’à 40% de notre borne supérieure !

En fait, même en répétant plusieurs fois des séries de 1000 ou 10.000 tests, nous ne sommes arrivés que rarement à produire un test qui trouve une valeur erronée avec le nombre d’itérations. D’où l’ajout du paramètre *CONFIDENCE_FACTOR*, qui réduit le nombre d’itérations effectuées par le programme.

Série de 10.000 tests (donc 400 erreurs attendues si on respecte la borne supérieure), avec un facteur de confiance de 10 :

```

1 confidence factor: 10
2 failures : 474 = 4.7% of tests
3 ...with an average of 1.22 absolute error
4 (Failed tests were removed from the following stats)
5 Three latest useful iterations (relative to total #iter) : 0.95, 0.93,
  0.85
6 Proportion of useful work : 0.15 (0.01 without confidence factor)

```

Cette fois, en faisant dix fois moins de calculs que prévus, on arrive quand même à respecter (ou presque) le taux de réussite qu’on s’était imposé. De plus, l’erreur commise n’est en moyenne pas très importante, ce qui pourrait être suffisant pour certaines applications nécessitant un calcul de coupe minimale.

On peut tout de même remarquer que, malgré notre facteur de confiance, 85% des itérations ont été calculées alors que la réponse avait déjà été trouvée. Ceci semble excessif mais est le prix à payer pour avoir une très grande majorité de résultats corrects (ceci peut se voir dans les trois itérations “gagnantes” les plus lointaines, qui sont toutes les trois arrivées vers la fin de l’algorithme).

Testons maintenant une valeur extrême de confiance (10.000 tests, facteur de confiance 100) :

```

1 confidence factor: 100
2 failures : 5139 = 51.39% of tests
3 ...with an average of 2.7 absolute error
4 Three latest useful iterations (relative to total #iter) : 0.5, 0.0,
  0.0
5 Proportion of useful work : 0.03 (0.0 without confidence factor)

```

(L’algorithme ne testant qu’une ou deux itérations par graphes, les deux dernières lignes moins significatives)

⁵En tout cas pour l’ensemble des graphes accessibles par notre fonction de génération aléatoire, qui n’est bien sûr qu’un sous-ensemble qu’on espère représentatif de tous les graphes possibles.

On a maintenant une majorité de résultats erronés (ce qui était attendu vu le facteur de confiance très élevé). Cependant, on peut quand même remarquer que l'erreur moyenne n'est pas si grande que ça (elle est significative, mais pas extravagante). On pourrait imaginer des applications nécessitant souvent d'avoir très rapidement une petite coupe d'un graphe, sans contrainte de minimalité, où l'algorithme de Karger avec un grand facteur de confiance serait parfaitement adapté.

4 Références

- [1] Karger, David (1993). "Global Min-cuts in RNC and Other Ramifications of a Simple Mincut Algorithm". Proc. 4th Annual ACM-SIAM Symposium on Discrete Algorithms.
- [2] "n points are chosen randomly and uniformly inside the unit square and pairs of points closer to each other than a predefined distance d are connected by an edge" <http://hal.elte.hu/~nepusz/development/igraph/tutorial/tutorial.html>

5 Annexes

5.1 Génération aléatoire de graphes

Pour avoir des résultats significatifs, il nous a semblé préférable de tester l'algorithme avec des graphes générés aléatoirement plutôt qu'avec des "cas particuliers" codés en dur dans le code.

L'algorithme utilisé, proposé par la bibliothèque `igraph`, prend comme paramètres :

n : le nombre de sommets du graphe

r : un nombre entre 0 et 1

Il génère le graphe en posant les n sommets aléatoirement dans un carré de côté 1, puis en reliant d'une arête les points situés à une distance inférieure à r .

L'algorithme peut donc générer des graphes non connexes, qui sont des cas triviaux jugés non intéressants car chaque itération de l'algorithme (modifié pour accepter ce cas) renverra la bonne valeur avec une probabilité de 1. Il a donc été choisi de détecter les graphes connexes renvoyés par l'algorithme et de demander un nouveau graphe aléatoire dans ce cas.

Au cours, on a calculé que l'algorithme de Karger avait une complexité en $O(n^4)$, il convient donc de choisir un n d'une taille raisonnable. Sur notre matériel et pour des séries de 1000 tests, 10 sommets semblent être une borne maximale pour avoir un temps de calcul viable.

Le choix de r , qui modifiera surtout le nombre d'arêtes du graphe, semblait avoir moins d'importance sur le temps de calcul. Mais il nous semblait préférable d'éviter les valeurs extrêmes (0 : aucunes arêtes, 1 : tous les sommets sont

reliés).

Au final, par plaisir de randomisation, il a été décidé de générer aléatoirement pour chaque graphe un n entre 5 et 10, et un r entre 0,4 et 1,0 (ces valeurs étant des constantes facilement modifiables).

5.2 Améliorations possibles

- Fonction de génération couvrant plus de types de graphes
- Graphique de la position relative des itérations “gagnantes” pour chaque test (information déjà présente dans les instances de *TestResult*, il reste “seulement” à l’afficher).
- En plus des valeurs moyennes, calculer les variances et les valeurs maximales des grandeurs mesurées.

5.3 Code source complet

Le code a normalement été fourni avec le présent rapport. Il est aussi disponible sur GitHub :

<https://github.com/BafAltom/info-f413-random-mincut>

```
1 #!/usr/bin/python
2 # -*- coding: utf-8 -*-
3
4 from igraph import *
5 # igraph: http://igraph.sourceforge.net/
6 import random
7 import math
8
9 random.seed()
10
11 # Nombre de graphes qui vont tre test s
12 NUMBER_OF_TESTS = 1000
13 # Proportion d'it rations qui ne vont pas tre effectues (par
14 # rapport la borne minimale vue au cours)
15 CONFIDENCE_FACTOR = 1 # def: 5
16 # Bornes min/max du nombre de sommets de chaque graphe
17 VERTICES.NB.MIN = 5 # def: 5
18 VERTICES.NB.MAX = 10 # def: 10
19 # Borne min/max du radius GRG de chaque graphe (une plus grande valeur
20 # signifie plus d'ar tes)
21 GRG.RADIUS.MIN = 0.4 # def: 0.4
22 GRG.RADIUS.MAX = 1.0 # def: 0.9
23 # Utilise (True) ou non (False) les graphes pour lesquels la coupe
24 # minimale n'a pas t trouve dans les statistiques
25 REMOVE_FAILED_TESTS.IN.STATS = True
26 # Affiche (True) ou non (False) un r sum des tests pour lesquels la
27 # coupe minimale n'a pas t trouve
28 DISPLAY_FAILURES = False
29 # Recommence l'ex cution si aucune faute n'a t trouve (peut
30 # prendre beaucoup de temps...)
31 CONTINUE.ON.PERFECT.TEST = False
32
33 class TestResult: # Classe stockant un r sum de l'ex cution d'un
34     test
35     def __init__(self, graph):
36         #self.graph = graph
37         self.number_of.vertices = graph.vcount()
```

```

32     self.number_of_edges = graph.ecount()
33     self.number_of_iterations = None
34     self.improving_iterations = []
35     self.improved_values = []
36     self.found_value = None
37     self.correct_value = None
38
39     def consistencyCheck(self):
40         assert(len(self.improving_iterations) == len(self.improved_values)
41                )
42         assert(self.found_value != None)
43         assert(self.correct_value != None)
44
45     def __str__(self):
46         self.consistencyCheck()
47         s = "Test with a graph of " + str(self.number_of_vertices) + "
48             vertices and " + str(self.number_of_edges) + " edges.\n"
49         s += "Used " + str(self.number_of_iterations) + " iterations.\n"
50         for i in range(len(self.improving_iterations)):
51             s += "Found value : " + str(self.improved_values[i]) + " at
52                 iteration " + str(self.improving_iterations[i]) + "\n"
53         s += "Final value : " + str(self.found_value) + "\n"
54         s += "Correct value : " + str(self.correct_value) + "\n"
55         return s
56
57 def random_mincut_one_instance(g):
58     while (g.vcount() > 2 and g.ecount() > 0):
59         edge_list = g.get_edgelist()
60         if (g.ecount() == 1): randedge_id = 0 # hotfix pour des graphes
61             non connexes
62         else: randedge_id = random.randint(0,g.ecount()-1)
63         chosen_edge = edge_list[randedge_id]
64         chosen_edge_source = chosen_edge[0]
65         chosen_edge_target = chosen_edge[1]
66         for e in edge_list:
67             if (e[0] == chosen_edge_target or e[1] == chosen_edge_target):
68                 if (e[0] == chosen_edge_target and e[1] !=
69                     chosen_edge_source):
70                     g.add_edge(chosen_edge_source, e[1])
71                 elif (e[1] == chosen_edge_target and e[0] !=
72                     chosen_edge_source):
73                     g.add_edge(e[0], chosen_edge_source)
74                 edge_id = g.get_eid(e[0], e[1])
75                 g.delete_edges(edge_id)
76
77         assert(g.degree(chosen_edge_target) == 0)
78         g.delete_vertices(chosen_edge_target)
79     return g.ecount()
80
81 def random_mincut_all(graph, myTestResult):
82     number_of_vertices = graph.vcount()
83     best_mincut_value = 999999
84     required_nbr_iter = int(number_of_vertices*(number_of_vertices-1)*
85                             math.log(number_of_vertices))
86     nbr_iter = max(1,required_nbr_iter/CONFIDENCE_FACTOR)
87     myTestResult.number_of_iterations = nbr_iter
88     for i in range(nbr_iter):
89         g_temp = graph.copy()
90         result = random_mincut_one_instance(g_temp)
91         if (result < best_mincut_value):
92             best_mincut_value = result
93         myTestResult.improving_iterations.append(i)

```

```

87         myTestResult.improved.values.append(best_mincut.value)
88
89     myTestResult.found_value = best_mincut.value
90     myTestResult.correct_value = int(graph.mincut().value)
91     return myTestResult
92
93 def generate_random_connex_graph(number_of_vertices, GRG_radius):
94     g = Graph.GRG(number_of_vertices, GRG_radius)
95     while (g.cohesion() == 0):
96         g = Graph.GRG(number_of_vertices, GRG_radius)
97     return g
98
99 def strRound2(aFloat):
100     return str(round(aFloat, 2))
101
102 if __name__ == "__main__":
103     anotherRound = True
104     while (anotherRound):
105         tests_results = []
106         for i in range(NUMBER_OF_TESTS):
107             if (NUMBER_OF_TESTS < 1000): print "test " + str(i+1) + " of "
108                 + str(NUMBER_OF_TESTS) + "..."
109             elif ((i+1)%(NUMBER_OF_TESTS/100)) == 0: print "test " + str(
110                 i+1) + " of " + str(NUMBER_OF_TESTS) + "..."
111
112             number_of_vertices = random.randint(VERTICES_NB_MIN,
113                 VERTICES_NB_MAX)
114             GRG_radius = random.randint(GRG_RADIUS_MIN*100, GRG_RADIUS_MAX
115                 *100)/float(100)
116             g = generate_random_connex_graph(number_of_vertices,
117                 GRG_radius)
118             result = TestResult(g)
119             random_mincut.all(g, result)
120             tests_results.append(result)
121
122             # compute some statistics on tests_results
123             current_prop.sum = 0
124             current_failure.sum = 0
125             current_failure.value = 0
126             current_max.useful = 0
127             current_max.useful_2 = 0
128             current_max.useful_3 = 0
129             for tr in tests_results:
130                 tr.consistencyCheck()
131                 #failure count and value
132                 if(tr.found_value != tr.correct_value):
133                     if (DISPLAY_FAILURES):
134                         print "_____ "
135                         print "Failure : "
136                         print tr
137                     current_failure.sum += 1
138                     current_failure.value += tr.found_value - tr.correct_value
139                     if (REMOVE_FAILED_TESTS_IN_STATS):
140                         continue
141                 # max useful
142                 number_of_correction = len(tr.improving.iterations)
143                 if (number_of_correction == 0):
144                     print tr
145                 last_useful.iteration_abs = tr.improving.iterations[
146                     number_of_correction-1]
147                 last_useful.iteration_rel = last_useful.iteration_abs / float(
148                     tr.number_of.iterations)

```



```

142         if (last.useful.iteration.rel > current_max.useful):
143             current_max.useful_3 = current_max.useful_2
144             current_max.useful_2 = current_max.useful
145             current_max.useful = last.useful.iteration.rel
146             # proportion of useful work
147             current_prop.sum += last.useful.iteration.rel
148
149     prop_average = None
150     if(REMOVE_FAILED.TESTS.IN.STATS):
151         prop_average = current_prop.sum/(len(tests.results) -
            current_failure.sum)
152     else:
153         prop_average = current_prop.sum/len(tests.results)
154
155     if (CONTINUE.ON.PERFECT.TEST and current_failure.sum == 0):
156         anotherRound = True
157     else: anotherRound = False
158
159     print "_____ "
160     print "confidence factor: " + str(CONFIDENCE_FACTOR)
161     print "failures : " + str(current_failure.sum) + " = " + strRound2
        (100*current_failure.sum/float(NUMBER_OF.TESTS)) + "% of tests"
162     if (current_failure.sum > 0): print "...with an average of " +
        strRound2(current_failure.value/float(current_failure.sum)) + "
        absolute error"
163     if (REMOVE_FAILED.TESTS.IN.STATS): print "(Failed tests were removed
        from the following stats)"
164     print "Three latest useful iterations (relative to total #iter) : "
        + strRound2(current_max.useful) + ", " + strRound2(
        current_max.useful_2) + ", " + strRound2(current_max.useful_3)
165     print "Proportion of useful work : " + strRound2(prop_average) + " (
        " + strRound2(prop_average/CONFIDENCE_FACTOR) + " without
        confidence factor)"

```