

IFT2015 automne 2013 — Devoir 1

Miklós Csűrös

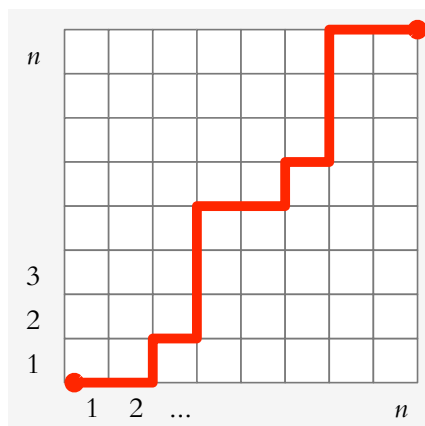
5 septembre 2013

À remettre avant 20 :15 le 12 septembre. Remettez un rapport écrit par email (à csuros@iro...) en format PDF. (La taille du fichier ne doit pas dépasser 2^{20} octets.)

Vous avez le droit de travailler en un équipe de 2 sur problème 3.d et le soumettre plus tard, avant 20 :15 le 16 septembre. Travaillez seul sur tous les autres problèmes.

1.1 Chemins sur une grille (15 points)

On veut compter les chemins sur une grille $n \times n$.



Un chemin sur la grille $n \times n$ est une séquence $(x_0, y_0), (x_1, y_1), \dots, (x_m, y_m)$ avec $(x_0, y_0) = (0, 0)$, $(x_m, y_m) = (n, n)$, et $(x_i, y_i) \in \{(x_{i-1} + 1, y_{i-1}), (x_{i-1}, y_{i-1} + 1)\}$ pour tout $i = 1, 2, \dots, m$. (Dans d'autres mots, on peut déplacer seulement vers le droit ou vers le haut par une case.)

a. (5 points) ► Démontrer que le nombre de chemins $C(n)$ est

$$C(n) = \binom{2n}{n} = \frac{(2n)!}{(n!) \cdot (n!)}.$$

b. (10 points) ► Utiliser la formule de Stirling pour caractériser la croissance de $\ln C(n)$ par une fonction f (aussi simple que possible) t. q. $\ln C(n) \sim f(n)$ ou $\lim_{n \rightarrow \infty} \frac{\ln C(n)}{f(n)} = 1$.

1.2 Comment encoder les entiers ? (35+15 points)



On veut transmettre une séquence de nombres naturels n_0, n_1, n_2, \dots ($n_i \geq 0$). La difficulté est qu'on ne peut utiliser que deux symboles, comme en code Morse : **0** et **1**. (Il s'agit de l'abstraction d'un fichier binaire.)

On définit donc un code $\rho(n)$ (séquence de symboles de **0** et **1**) pour tout n , et on transmet les codes $\rho(n_0), \rho(n_1), \dots$, l'un après l'autre. En plus d'être *déchiffrable* (on peut récupérer la séquence n_0, n_1, \dots sans ambiguïté à partir de la chaîne concaténée), ρ devrait être *instantané* : on veut déchiffrer n_i dès qu'on reçoit tous les bits de $\rho(n_i)$.

Encodage unaire. L'exemple le plus simple est le *code unaire* $u(n)$:

$$u(n) = \underbrace{1\ 1\ \dots\ 1\ 0}_{n\ \text{fois}} \quad (1.1)$$

Clairement, il est possible de décoder une séquence de codes concaténés $u(n_0), u(n_1), \dots$, car il suffit de compter les **1**s jusqu'au premier **0**. Mais ce code n'est pas économique de tout ...

Encodage binaire. L'encodage binaire standard $b(n)$ est beaucoup plus efficace, parce qu'il utilise seulement d symboles où d est le plus petit nombre naturel tel que $n < 2^d$. Pour $n = \sum_{i=0}^{d-1} x_i \cdot 2^i$ (avec $x_i = 0, 1$) on transmet la séquence $x_{d-1}, x_{d-2}, \dots, x_0$, en commençant par le bit de poids fort. En particulier $b(0) = \lambda$ (chaîne vide), $b(1) = 1$, $b(2) = 1\ 0$, etc.

Par contre, b n'est pas déchiffrable : $(1, 2, 1)$ et $(6, 1)$ mènent à la même séquence **1 1 0 1**.

Encodage f. On peut construire un code instantané en préfixant $b(n)$ par sa longueur encodé en unaire : si $b(n)$ comprend d symboles, alors on écrit d fois **1**, suivi par un seul **0**, et l'encodage binaire sur d bits.

$$f(n) = u(|b(n)|) \oplus b(n), \quad (1.2)$$

où \oplus dénote la concaténation et $|\dots|$ la longueur en bits. Le code f est instantanément déchiffrable : compter les **1**s au début ($= d$), jusqu'au premier **0**, et décoder les d bits suivants contenant $b(n_i)$.

Encodage g. On peut améliorer le code f en spécifiant le nombre de bits avec un encodage plus efficace que l'unaire. On se sert de la récursivité :

$$g_0(n) = \begin{cases} \lambda & \text{chaîne de longueur 0 quand } n = 0 \\ g_0(|b(n)| - 1) \oplus b(n) & \{n > 0\} \end{cases} \quad (1.3)$$

On rend le code déchiffrable par un bit **0** apposé : on transmet $g_1(n) = g_0(n) \oplus 0$. Par exemple, **1 1 0 1 1 1 0 1 1 0 0 0** encode la séquence 7, 2, 0.

n	$b(n)$	$f(n)$	$g_0(n)$
0	λ	0	λ
1	1	1 0 1	1
2	1 0	1 1 0 1 0	1 1 0
3	1 1	1 1 0 1 1	1 1 1
4	1 0 0	1 1 1 0 1 0 0	1 1 0 1 0 0
5	1 0 1	1 1 1 0 1 0 1	1 1 0 1 0 1
6	1 1 0	1 1 1 0 1 1 0	1 1 0 1 1 0
7	1 1 1	1 1 1 0 1 1 1	1 1 0 1 1 1
8	1 0 0 0	1 1 1 1 0 1 0 0 0	1 1 1 1 0 0 0 0
16	1 0 0 0 0 0	1 1 1 1 1 0 1 0 0 0 0 0	1 1 0 1 0 0 1 0 0 0 0 0
2015	1 1 1 1 1 0 1 1 1 1 1	[23 symboles] 1 1 1 1 0 1 0 1 1 1 1 1 0 1 1 1 1 1	1 1 1 1 1 0 1 1 1 1 1

a. (9 points) ► Combien de symboles doit on utiliser pour écrire 10^{100} (un *googol*) dans les encodages différents u , b , f et g_1 ? ► En général, quelle est la longueur exacte de u , b , et f en fonction de n ?

b. (13 points) ► Donner un algorithme *récurif* $encodeBeta(n)$ qui affiche l'encodage b d'un nombre naturel n .

c. (13 points) ► Donner un algorithme *récurif* $encodeOmega(n)$ qui émet l'encodage $g_1(n)$. (Vous avez le droit d'appeler $encodeBeta$ de **b**.)

d. (15[♥] points boni) Étudiez l'efficacité de la méthode de compression suivante qu'on va appeler *intZip*. Prendre un fichier spécifié à l'entrée, et compter la fréquence de chaque symbole¹ dedans. Encoder le symbole le plus fréquent par $g_1(0)$, le deuxième plus fréquent par $g_1(1)$, etc. : le k -ème plus fréquent par $g_1(k - 1)$ en général. Le fichier à la sortie est la séquence de g_1 pour les symboles de l'entrée.

Dans votre étude, prenez quelques fichiers «typiques» dans un contexte quelconque de taille $\geq 1\text{Mb}$. Vous pouvez vous servir de la collection de données à <http://introcs.cs.princeton.edu/java/data/> ou n'importe quel autre jeu de données disponible sur l'Internet. Comparez la longueur de compression *intZip* avec la taille du fichier comprimé par outils populaires (p.e., *gzip*, *bzip*) de votre choix. Soumettez votre code pour calculer la longueur de compression (notez qu'il n'est pas nécessaire de construire le code g) en un seul archive JAR, et un rapport de vos expériences en un fichier PDF.

REMARQUE. Lire le fichier deux fois comme ici n'est pas pratique. Les méthodes de compression qui utilisent l'encodage f ou g les appliquent à une séquence d'entiers produites «à la volée», en lisant le fichier de l'entrée. La clé est de choisir un heuristique qui transforme la séquence d'entrée en entiers typiquement petits pour exploiter les propriétés de g .

¹Choisissez ce que c'est un «symbole» selon le contexte. Les caractères (**char**) ou les octets (**byte**) du fichier peuvent suffire, mais parfois (ADN, musique, etc.) c'est mieux de définir un symbole comme un bloc de caractères consécutifs.