

PETER KIRSCHENHOFER

HELMUT PRODINGER

Approximate counting : an alternative approach

Informatique théorique et applications, tome 25, n° 1 (1991), p. 43-48

<http://www.numdam.org/item?id=ITA_1991__25_1_43_0>

© AFCET, 1991, tous droits réservés.

L'accès aux archives de la revue « Informatique théorique et applications » implique l'accord avec les conditions générales d'utilisation (<http://www.numdam.org/conditions>). Toute utilisation commerciale ou impression systématique est constitutive d'une infraction pénale. Toute copie ou impression de ce fichier doit contenir la présente mention de copyright.

NUMDAM

Article numérisé dans le cadre du programme
Numérisation de documents anciens mathématiques

<http://www.numdam.org/>

APPROXIMATE COUNTING: AN ALTERNATIVE APPROACH (*)

by Peter KIRSCHENHOFER ⁽¹⁾ and Helmut PRODINGER ⁽¹⁾

Abstract. – In this note an alternative analysis of approximate counting is presented by using a lemma from the calculus of finite differences instead of the Mellin integral transform.

Résumé. – Cette note présente une analyse de l'algorithme de comptage approximatif qui repose sur un lemme du calcul des différences finies au lieu d'une utilisation de la transformée de Mellin.

R. Morris [6] has proposed a probabilistic algorithm that maintains an approximate count in the interval 1 to n using only about $\log_2 \log_2 n$ bits. Approximate counting (with a binary base) starts with counter $C=1$. After n increments C should contain a good approximation to $\lfloor \log_2 n \rfloor$; thus C should be increased by 1 after another n increments approximately. Since only C is known the algorithm has to base its decision on the content of C alone. The following principle is used to increment the counter:

$$C := C + \begin{cases} 0 & \text{with probability } 1 - 2^{-C} \\ 1 & \text{with probability } 2^{-C} \end{cases}$$

Compare Flajolet [2] for a more detailed description. In the same paper a detailed analysis of this algorithm is presented and the following results are shown:

THEOREM 1: *After n successive increments the average content \bar{C}_n of the counter satisfies:*

$$\bar{C}_n \sim \log_2 n + \frac{\gamma}{\log 2} - \alpha + \frac{1}{2} + \delta_1(\log_2 n);$$

(*) Received January 1989, accepted in January 1990.

We thank an anonymous referee for his/her guidance concerning a better presentation of this paper.

⁽¹⁾ Department of Algebra and Discrete Mathematics Technical University of Vienna, Austria.

the variance σ_n^2 satisfies

$$\sigma_n^2 \sim \sigma_\infty^2 + \delta_2(\log_2 n)$$

where

$$\sigma_\infty^2 = \frac{\pi^2}{6 \log^2 2} - \alpha - \beta + \frac{1}{12} - \frac{1}{\log 2} \sum_{k \geq 1} \frac{1}{k \sinh(\theta k)}$$

with

$$\alpha = \sum_{k \geq 1} \frac{1}{2^k - 1}, \quad \beta = \sum_{k \geq 1} \frac{1}{(2^k - 1)^2}, \quad \theta = \frac{2\pi^2}{\log 2},$$

and $\delta_1(x)$ as well as $\delta_2(x)$ are periodic functions with period 1, mean value 0, and amplitude less than 10^{-4} .

Flajolet achieves this result by computing the probability $p_{n,l}$ of having counter value l after n increments:

$$p_{n,l} = \sum_{j=0}^{l-1} \frac{(-1)^j 2^{-\binom{j}{2}} (1 - 2^{-(l-j)})^n}{Q_j Q_{l-1-j}}, \quad (1)$$

where

$$Q_n = \left(1 - \frac{1}{2}\right) \left(1 - \frac{1}{4}\right) \dots \left(1 - \frac{1}{2^n}\right). \quad (2)$$

The quantities

$$\bar{C}_n = \sum_{l=1}^n l p_{n,l} \quad \text{and} \quad \sigma_n^2 = \sum_{l=1}^n l^2 p_{n,l} - (\bar{C}_n)^2 \quad (3)$$

are then analyzed using real analysis and certain properties of the Mellin integral transform.

In the sequel we analyze the quantities in (3) by an alternative approach which avoids completely the unpleasant real analysis by making use of a classical lemma from the calculus of finite differences. This method has been attributed by Knuth [5] to S. O. Rice.

LEMMA 2: Let C be a curve surrounding the points $1, 2, \dots, n$ in the complex plane and let $f(z)$ be analytic inside C . Then

$$\sum_{k=1}^n \binom{n}{k} (-1)^k f(k) = -\frac{1}{2\pi i} \int_C \langle n; z \rangle f(z) dz,$$

where

$$\langle n; z \rangle = \frac{(-1)^{n-1} n!}{z(z-1)\dots(z-n)}.$$

Moving the contour of integration it turns out that the asymptotic expansion of the alternating sum is obtained via

$$\sum \text{Res}(\langle n; z \rangle f(z)) \tag{4}$$

where the sum is taken over all poles different from $1, \dots, n$. (Flajolet and Sedgewick have used this technique to simplify the analysis of digital search trees in [3].)

In order to apply Lemma 2 we write (1) as

$$p_{n,l} = \sum_{k=0}^n \binom{n}{k} (-1)^k \sum_{j=0}^{l-1} \frac{(-1)^j 2^{-\binom{j}{2}}}{Q_j Q_{l-1-j}} 2^{-(l-j)k}. \tag{5}$$

By Euler's *partition identities* [1]

$$\sum_{j \geq 0} \frac{(-1)^j 2^{-\binom{j}{2}} x^j}{Q_j} = \prod_{m \geq 0} \left(1 - \frac{x}{2^m}\right)$$

and

$$\sum_{j \geq 0} \frac{x^j}{Q_j} = \prod_{m \geq 0} \left(1 - \frac{x}{2^m}\right)^{-1},$$

the inner sum in (5) is the coefficient

$$[x^{l-1}] \sum_{j \geq 0} \frac{(-1)^j 2^{-\binom{j}{2}} x^j}{Q_j} \sum_{s \geq 0} 2^{-(s+1)k} \frac{x^s}{Q_s} = 2^{-k} [x^{l-1}] \prod_{m=0}^{k-1} \left(1 - \frac{x}{2^m}\right).$$

Thus

$$\bar{C}_n = \sum_{l=1}^n l p_{n,l} = \sum_{k=0}^n \binom{n}{k} (-1)^k 2^{-k} \sum_{l=1}^n l [x^{l-1}] \prod_{m=0}^{k-1} \left(1 - \frac{x}{2^m}\right),$$

where the summation on l can be extended to infinity without changing the value. Since

$$\sum_{l \geq 1} l [x^{l-1}] g(x) = g'(1) + g(1)$$

we get immediately the following identity for \bar{C}_n which we state as a Proposition, since it is interesting in its own right.

PROPOSITION 3

$$\bar{C}_n = 1 - \sum_{k=0}^n \binom{n}{k} (-1)^k 2^{-k} Q_{k-1}. \tag{6}$$

With

$$Q_z = \frac{\prod_{j \geq 1} (1 - (1/2^j))}{\prod_{j \geq 1} (1 - (1/2^{z+j}))}$$

we may write (6) as

$$\bar{C}_n = 1 - \sum_{k=1}^n \binom{n}{k} (-1)^k f(k) \quad \text{with } f(z) = 2^{-z} Q_{z-1}. \tag{7}$$

The most significant pole of $\langle n; z \rangle f(z)$ different from $1, \dots, n$ is $z=0$, and the residue is

$$-\frac{H_n}{L} + \frac{1}{2} + \alpha \sim -\log_2 n - \frac{\gamma}{L} + \frac{1}{2} + \alpha,$$

with $L = \log 2$ and H_n the n -th Harmonic number. Further poles in $z = 2k\pi i/L$ with residue equal to

$$\frac{1}{L} \Gamma\left(\frac{2k\pi i}{L}\right)$$

give rise to the periodic fluctuation $\delta_1(x)$.

The next significant poles have real part $z = -1$ (simple!), so that the error term is of order n^{-1} ,

A similar approach allows us to evaluate the variance σ_n^2 in a few lines. Since

$$\sum_{l \geq 1} l^2 [x^{l-1}] g(x) = g^n(1) + 3g'(1) + g(1)$$

we have

$$\sum_{k=1}^n l^2 p_{n,l} = 1 + \sum_{k=1}^n \binom{n}{k} (-1)^k f(k)$$

with

$$f(k) = 2^{-k} Q_{k-1} (2 T_{k-1} - 3), \quad T_{k-1} = \sum_{i=1}^{k-1} \frac{1}{2^{i-1}}.$$

The continuation of Q_k has been referred above; T_k can be continued via

$$T_z = \alpha - \sum_{i \geq 1} \frac{1}{2^{z+i-1}}.$$

The most important pole is again $z=0$ (triple) with residue

$$\begin{aligned} & \frac{2}{L^2} \left(\frac{H_n^2 + H_n^{(2)}}{2} + H_n L \left(\frac{1}{2} - \alpha \right) \right) + \alpha^2 - 4\alpha - \beta - \frac{2}{3} \\ & \sim \log_2^2 n + (\log_2 n) \left(\frac{2\gamma}{L} - 2\alpha + 1 \right) + \alpha^2 - 2\alpha - \beta + \frac{\gamma}{L} + \frac{\gamma^2}{L^2} - \frac{2\alpha\gamma}{L} + \frac{\pi^2}{6L^2} - \frac{2}{3}. \end{aligned}$$

Subtracting \bar{C}_n^2 yields for the variance (apart from smaller order terms)

$$\frac{\pi^2}{6L^2} - \alpha - \beta + \frac{1}{12} - \frac{2}{L^2} \sum_{k \geq 1} \left| \Gamma \left(\frac{2k\pi i}{L} \right) \right|^2 + \delta_2(\log_2 n),$$

where the series originates from the mean of $\delta_1^2(x)$. Observe that

$$|\Gamma(iy)|^2 = \frac{\pi}{y \sinh \pi y}, \quad y \in \mathbf{R},$$

to obtain the form of σ_∞^2 , whence Theorem 1.

Finally, it is striking to note that the constant σ_∞^2 , whose numerical value is

$$\sigma_\infty^2 = 0.7630141871107195182 \dots$$

differs from the simpler

$$\frac{1}{2 \log 2} + \frac{1}{24} = 0.76301418711114837034 \dots$$

only in the twelfth digit after the decimal point.

Indeed, using a transformation result for Dedekind's η -function

$$\eta(\tau) = e^{\pi i \tau / 12} \prod_{n \geq 1} (1 - e^{2\pi i n \tau}), \quad \Im \tau > 0,$$

the constant σ_{∞}^2 appearing in (8) turns out to be

$$\frac{1}{2L} + \frac{1}{24} - h_1\left(\frac{2\pi^2}{L}\right) + \frac{2\pi^2}{L^2} h_2\left(\frac{4\pi^2}{L}\right),$$

where

$$h_1(x) = \sum_{k \geq 1} \frac{(-1)^{k-1}}{k(e^{kx} - 1)} \quad \text{and} \quad h_2(x) = \sum_{k \geq 1} \frac{e^{kx}}{(e^{kx} - 1)^2},$$

so that the numerical coincidence follows by the smallness of the involved functions $h_1(x)$ and $h_2(x)$.

Compare [4], where a closely related constant with relevance to digital search trees was evaluated by the authors and J. Schoissengeier.

Finally it should be noted that the expansion to arbitrary bases of the counter can be analyzed in the same style.

REFERENCES

1. G. E. ANDREWS, The Theory of Partitions, *Addison Wesley*, 1976.
2. P. FLAJOLET, Approximate Counting: A detailed Analysis, *BIT*, 1985, 25, pp. 113-134.
3. P. FLAJOLET and R. SEDGEWICK, Digital Search Trees Revisited, *S.I.A.M. J. Comput.*, 1986, 15, pp. 748-767.
4. P. KIRSCHENHOFER, H. PRODINGER and J. SCHOISSENGEIER, Zur Auswertung gewisser numerischer Reihen mit Hilfe modularer Funktionen, in *Zahlentheoretische Analysis II*, E. HLAJKA ed., Springer, Berlin, 1987, pp. 108-110.
5. D. E. KNUTH, The Art of Computer Programming, 3, *Addison Wesley*, 1973.
6. R. MORRIS, Counting Large Numbers of Events in Small Registers, *Comm. A.C.M.*, 1978, 21, pp. 840-842.