

A unified approach to linear probing hashing

Svante Janson

Uppsala University

Alfredo Viola Universidad de la República

FoCM 2014 Dedicated to Philippe Flajolet

Bucket parking scheme

- Each parking lots can hold up to r cars
- Related to analysis of bucket hashing algorithms



Further research: Bucket parking scheme

Bucket parking scheme

- Each parking lots can hold up to r cars
- Related to analysis of bucket hashing algorithms



Further research: Bucket parking scheme

Bucket parking scheme

- Each parking lots can hold up to r cars
- Related to analysis of bucket hashing algorithms



Further research: Bucket parking scheme

Bucket parking scheme

- Each parking lots can hold up to r cars
- Related to analysis of bucket hashing algorithms



Further research: Bucket parking scheme

Bucket parking scheme

- Each parking lots can hold up to r cars
- Related to analysis of bucket hashing algorithms



Further research: Bucket parking scheme

Bucket parking scheme

- Each parking lots can hold up to r cars
- Related to analysis of bucket hashing algorithms



Bucket parking scheme

- Each parking lots can hold up to r cars
- Related to analysis of bucket hashing algorithms



Bucket parking scheme

- Each parking lots can hold up to r cars
- Related to analysis of bucket hashing algorithms



Bucket parking scheme

- Each parking lots can hold up to r cars
- Related to analysis of bucket hashing algorithms



Further research: Bucket parking scheme

Bucket parking scheme

- Each parking lots can hold up to r cars
- Related to analysis of bucket hashing algorithms



The problem.

- The table has m places of capacity b to hash (park) from 0 to m 1, and n inserted elements (cars).
- Each element is given a hash value (preferred parking lot).
- If place is not full, then the element is stored there.
- Otherwise, looks sequentially for an empty place.
- If no empty place up to the end of the table, the search follows at location 0.
- Several R.V. to study, mainly related with cost of individual searches and total construction cost.
- Very important special case: Parking Problem.
- In parking the car is lost if no available place to park up to the end of the table.
- Main R.V. is the number of lost cars.

Linear Probing Hashing.



MATHÉMATIQUES ET INFORMATIQUE

Hachage, arbres, chemins & graphes Philippe Chassaing[†] et Philippe Flajolet[‡]

The mathematical beauty of Linear Probing!

Mathématiques discrètes et continues se rencontrent et se complètent volontiers harmonieusement. C'est cette thèse que nous voudrions illustrer en discutant un problème classique aux ramifications nombreusesl'analyse du hachage avec essais linéaires. L'exemple est issu de l'analyse d'algorithmes, domaine fondé par Knuth et qui se situe lui-même « à cheval » entre l'informatique, l'analyse combinatoire, et la théorie des probabilités. Lors de son traitement se croisent au fil du temps des approches très diverses, et l'on rencontrera des questions posées par Ramanujan à Hardy en 1913, un travail d'été de Knuth datant de 1962 et qui est à l'origine de l'analyse d'algorithmes en informatique, des recherches en analyse combinatoire du statisticien Kreweras, diverses rencontres avec les modèles de graphes aléatoires au sens d'Erdös et Rényi, un peu d'analyse complexe et d'analyse asymptotique, des arbres qu'on peut voir comme issus de processus de Galton-Watson particuliers, et, pour finir, un peu de processus, dont l'ineffable mouvement Brownien! Tout ceci contribuant in fine à une compréhension très précise d'un modèle simple d'aléa discret.

1962: Summer work by Don Knuth ...







NOTES ON "OPEN" ADDRESSING.

D. Knuth 7/22/63

8354

1. Introduction and Definitions. Upon addressing is a widely-used technique for keeping "symbol tables." The method was first used in 1954 by Samuel, Andahl. and Boohne in an assembly program for the IBM TOL. An extensive discussion of the method was given by Peterson in 1997 [1], and frequent references have been made to it ever since (e.g. Schay and Spruth [2], Iverson [3]). However, the timing characteristics have apparently never been exactly established, and indeed the author has heard reports of several reputable mathematicians who failed to find the solution after some trial. Therefore it is the purpose of this note to indicate one way by which the solution can be obtained.

By first analysis of an algorithm, my maly a during summer 1962 in Medison .

We will use the following abstract model to describe the asthod: 5 is a positive integer, and we have an array of : variables x1,x2,...,xg. At the beginning, x; = 0, for 15158.

To "enter the k-th item in the table," we rean that an integer a. is calculated, 1 \$ 4, \$ N, depending only on the item, and the following process is carried out:

1. Set j = $a_{\underline{k}}$. 2. "The comparison step." I: $x_{\underline{j}}$ = 0, set $x_{\underline{j}}$ = 1 and stop; we say "the k-th item has fallen into position x1.

3. If j = H, go to step 5.

Increase j by 1 and return to step 2.
 "The overflow step." If this step is entered twice, the table is full,

i.e. x, = 1 for 1 \$ 1 \$ N. Otherwise set j to 1 and return to step 2.

Observe the cyclic character of this algorithm.

We are concerned with the statistics of this method, with respect to the number of times the comparison step must be executed. More precisely, we cansider all of the M" possible sequences a, eq...a, to be equally probable, and we ask, "What is the probability that the comparison step is used precisely m times when the k-th itan is placed?

2. Non-overflow (self-contained) sequences.

Let [k] denote the number of sequences $a_1, a_2, \dots a_k$ (1 $\leq a_1 \leq n$) in which no overflow step occurs during the entire process of placing k items, if the algorithm is used for N = n. (By convention, we net [1] = 1.)

erms 1: If
$$a_k \ge n + 1$$
, then $\begin{bmatrix} n \\ k \end{bmatrix} = (n+1)^k - |c(n+1)^{k-1}|$

Proof: This proof is based on the fact that $\begin{bmatrix} a \\ a \end{bmatrix}$ is precisely the number of sequences b_1, b_2, \dots, b_k (1 $\leq b_1 \leq \dots + 1$) in which, if the algorithm is carried out for $\beta = \alpha + 1$, then $x_{n+1} = 0$ at the end of the operation. This follows because every sequence of the former type is one of the latter, and conversely the condition in the intervalue interval $\beta \ge 0$, and that no overflow step occurs.

But sequences of the latter type are easily enumerated, because the algorithm has circular symmetry; of the $(n+1)^n$ possible sequences b_1,b_2,\ldots,b_k , exactly k/(n+1) of these lwave $x_{n+1}\neq 0$. This shows that

 $\begin{bmatrix} n \\ k \end{bmatrix} = (n+1)^k \left(1 - \frac{k}{k}\right).$

Original results.

- Let a hash table with m positions and n inserted elements.
- Let $P_{m,n}$ the probability of the last position being empty. • $P_{m,n} = (1 - \frac{n}{2})$

$$P_{m,n}=\left(1-\frac{n}{m}\right).$$

• Let $C_{m,n}$ the R.V. for the number of successful searches of a random element.

•
$$E[C_{m,n}] = \frac{1}{2} (1 + Q_0((m, n-1)))$$
.

•
$$E[C_{m,\alpha m}] = \frac{1}{2} \left(1 + \frac{1}{1-\alpha}\right)$$
 with $0 \le \alpha < 1$.

- $E[C_{n,n}] = \sqrt{\frac{\pi n}{8}} + O(1)$ (proved on may 20, 1965).
- Let $U_{m,n}$ the R.V. for the number of unsuccessful searches of a random element.

•
$$E[U_{m,n}] = \frac{1}{2} (1 + Q_1((m, n - 1))).$$

• $E[U_{m,\alpha m}] = \frac{1}{2} \left(1 + \frac{1}{(1-\alpha)^2}\right)$ with $0 \le \alpha < 1$.

ullet The Ramanujan Q function is the special case $Q_0(n,n)$ of

$$Q_r(m,n) = \sum_{k=0}^n inom{k+r}{k} rac{n^{\underline{k}}}{m^n}$$

Linear Probing and the Symbolic Method (I).



Linear Probing and Graphs

Donald E. Knuth, Stanford University

Dedicated to Philippe Patrick Michel Flajolet

Abstract. Mallows and Riordan showed in 1968 that labeled trees with a small number of inversions are related to labeled graphs that are connected and sparse. Wright enumerated sparse connected graphs in 1977, and Kreweras related the inversions of trees to the so-called "parking problem" in 1980. A combination of these three results leads to a surprisingly simple analysis of the behavior of hashing by linear probing, including higher moments of the cost of successful search. The purpose of this note is to exhibit a surprisingly simple solution to a problem that appears in a recent book by Sedgewick and Flajolet [9]:

Exercise 8.39 Use the symbolic method to derive the EGF of the number of probes required by linear probing in a successful search, for fixed M.

The authors admitted that they did not know how to solve the problem, in spite of the fact that a "symbolic method" was the key to the analysis of all the other algorithms in their book. Indeed, the second moment of the distribution of successful search by linear probing was unknown when [9] was published in 1996.

Conclusions (Knuth).

7. Personal remarks. The problem of linear probing is near and dear to my heart, because I found it immensely satisfying to deduce (5.4) when I first studied the problem in 1962. Linear probing was the first algorithm that I was able to analyze successfully, and the experience had a significant effect on my future career as a computer scientist. None of the methods available in 1962 were powerful enough to deduce the expected square displacement, much less the higher moments, so it is an even greater pleasure to be able to derive such results today from other work that has enriched the field of combinatorial mathematics during a period of 35 years.

The reader will note that Sedgewick and Flajolet's exercise 8.39 has not truly been solved, strictly speaking, because we have not found the EGF $\sum_{n=0}^{m-1} F_{mn}(x) z^n/n!$ as requested. However, Sedgewick and Flajolet should be happy with any analysis of linear probing that uses symbolic methods associated with generating functions in an informative way.

Linear Probing and the Symbolic Method (II).







Algorithmica (1998) 22: 490-515



On the Analysis of Linear Probing Hashing¹

P. Flajolet,2 P. Poblete,3 and A. Viola4

Combinatorial interpretation.

29	69	10		24		36	77	18	58
49	79					56	97	38	78
0	1	2	3	4	5	6	7	8	9

 Any Linear Probing Hash table can be seen as a sequence of almost full tables (a subtable with all but the last bucket full).

• **Example:** [3-3],[4-4],[5-5],[6-2].

• This interpretation can be nicely handled by Analytic Combinatorics, since for example, it implies that it is enough to study almost full tables, and then use the *sequence* construction.



ON THE ANALYSIS OF LINEAR PROBING HASHING



Philippe Flajolet, INRIA Rocquencourt (France)



"My first analysis of an algorithm originally done during Summer 1962 at Madison."





$$\begin{array}{rcl} \textbf{CONSTRUCTIONS} \\ \textbf{Dictionary (I)} \\ \mathcal{F} & \mapsto & \{f_n\} & \mapsto & f(z) = \sum_n f_n \frac{z^n}{n!} \\ \\ \frac{1}{1-f} = 1+f+f^2+f^3+\cdots \\ exp(f) = 1+f+\frac{1}{2!}f^2+\frac{1}{3!}f^3+\cdots \\ exp(f) = 1+f+\frac{1}{2!}f^2+\frac{1}{3!}f^3+\cdots \\ \textbf{A} \cup \textbf{B} & \mapsto & A(z) + B(z) \\ \textbf{A} \times \textbf{B} & \mapsto & A(z) \times B(z) \\ \textbf{Seq A} & \mapsto & \frac{1}{1-A(z)} \\ \textbf{Set A} & \mapsto & exp(A(z)) \\ \textbf{Cycle A} & \mapsto & \log \frac{1}{1-A(z)} \end{array}$$

L.P.H.: Generating functions 4

Almost-full tables n = m - 1 have tree decomposition.



A nonlinear ODE translates Linear Probing

Table = Table '*' Table

$$F$$
 = $\int (zF)'F$
 $T'(1-\frac{1}{T})$ = $\frac{1}{z}$ with $T = zF$

Lemma 1.

$$F(z) = \frac{1}{z}T(z)$$
 where $T = ze^{T}$

Lemma 2. [Lagrange + Eisenstein + Cayley]

$$T(z) = z + 2\frac{z^2}{2!} + 9\frac{z^3}{3!} + \dots = \sum_n n^{n-1}\frac{z^n}{n!}$$

Defines the Tree function T(z)

Tree = Root * Set(Tree)

[Knuth63] The number of almost full tables (n keys)

$$F_n = (n+1)^{(n-1)}$$





EXAMPLE. [Euler] $(\exp(z))^{-1} = (\exp(-z))$

$\left(\sum \frac{z^n}{n!}\right)^-$	$^{1} = \left(\sum \frac{\left(-z\right)^{n}}{n!}\right)$
$\left(\sum \frac{z^n}{[n]!}\right)^{-1} =$	$\left(\sum q^{n(n-1)/2} \frac{(-z)^n}{[n]!}\right)$

$$\begin{split} \partial_z F(z) &= \partial_z (zF(z)) \cdot F(z) \\ \partial_z F(z,q) &= \left(\frac{F(z,q) - qF(qz,q)}{1-q}\right) \cdot F(z,q) \\ \partial_z F &= \mathrm{H} \ F \cdot F \end{split}$$

Moments result from

— ∂_q , differentiation w.r.t. q

— U, setting q = 1

For the r-th moment, apply $U \partial_q^r$ to difference differential equation

<u>Lemma</u>. Commutation rule with $Z[f] = z \cdot f$

 $\begin{array}{rcl} \mathbf{U} \mathbf{H} &=& \partial_z \mathbf{Z} \\ \mathbf{U} \, \partial_q \, \mathbf{H} &=& \mathbf{U} \, \partial_z \mathbf{Z} \partial_q + \frac{1}{2} \mathbf{Z} \partial_z^2 \mathbf{Z}, \\ \mathbf{U} \, \partial_q^2 \, \mathbf{H} &=& \partial_z \mathbf{Z} \, \mathbf{U} \, \partial_q^2 + \mathbf{Z} \partial_z^2 \mathbf{Z} \, \mathbf{U} \, \partial_q + \frac{1}{3} \mathbf{Z}^2 \partial_z^3 \mathbf{Z} \, \mathbf{U} \end{array}$

Proof. $\partial_q[n+1]$, Leibniz, &c.

• Let F_{bi+d} be the number of ways to construct an almost full table of length i+1 and size bi+d with $0 \le d \le b-1$. Then,

$$F_d(u):=\sum_{i\geq 0}F_{bi+d}rac{u^{bi+d}}{(bi+d)!}, \ \ \ N_d(z,w):=\sum_{s=0}^{b-1-d}w^{b-s}F_s(zw).$$

• $N_d(z, w)$ is the generating function for the number of almost full tables with more than d empty locations in the last bucket [Blake and Konheim 1977], [V. 2010].

$$\sum_{d=0}^{b-1}F_d(bz)x^d=x^b-\prod_{j=0}^{b-1}\left(x-rac{T(\omega^jz)}{z}
ight),
onumber \ N_d(bz,w)=\left[x^d
ight]rac{\prod\limits_{j=0}^{b-1}\left(1-xrac{T(\omega^jzw)}{z}
ight)-\prod\limits_{j=0}^{b-1}\left(1-rac{T(\omega^jzw)}{z}
ight)}{1-x},$$

where T is the Tree function and ω is a b-th root of unity.

- Let Q_{m,n,d} be the number of ways of inserting n elements into m buckets, where the last bucket has more than d empty slots.
- By a direct application of the sequence construction we have

$$\Lambda_d(bz,w):=\sum_{m\geq 0}\sum_{n\geq 0}Q_{m,n,d}rac{(bz)^n}{n!}w^{bm}=1+rac{N_d(bz,w)}{1-N_0(bz,w)}.$$

 Λ₀(z, w) is the generating function for the number of ways to construct hash tables such that their last bucket is not full.



• Let P be a property (e.g. cost of a successful search or block length).



• Let $p_{bi+d}(q)$ be the probability generating function of P calculated in the last cluster. Then

$$p_{m,n}(q) = \sum_{d=0}^{b-1} \sum_{i \ge 0} \binom{n}{bi+d} Q_{m-i-1,n-bi-d,0} \, F_{bi+d} \, p_{bi+d}(q).$$

• As a consequence,

$$egin{aligned} P(z,w,q) &:= \sum\limits_{m,n \geq 0} p_{m,n}(q) \; w^{bm} rac{z^n}{n!} = rac{\hat{N}_0(z,w,q)}{1-N_0(z,w)}, \qquad ext{with} \ \hat{N}_0(z,w,q) &:= \sum\limits_{d=0}^{b-1} w^{b-d} \; \sum\limits_{i \geq 0} F_{bi+d} \; rac{(zw)^{bi+d}}{(bi+d)!} \; p_{bi+d}(q), \end{aligned}$$

which could be directly derived with the sequence construction.

•
$$\hat{N}_0(z,w,1) = N_0(z,w), \ P(z,w,1) = \Lambda_0(z,w) - 1.$$

• Moreover, by singularity analysis of $\Lambda_0(lpha,y^{1/b}e^{-lpha},q)$ at y= 1,

$$\lim_{m\to\infty}\mathbf{P}_m[Q_{m,n,0}/m^n;b\alpha]=T_0(b\alpha)=\frac{b(1-\alpha)}{\prod_{j=1}^{b-1}\left(1-\frac{T(\omega^j\alpha e^{-\alpha})}{\alpha}\right)}.$$

• As a consequence,

$$\lim_{m\to\infty}\mathbf{P}_m[p_{m,n}(q)/m^n;b\alpha] \hspace{2mm} = \hspace{2mm} T_0(b\alpha)\tilde{N}_0(b\alpha,e^{-\alpha},q).$$

Summary of the q-calculus.

	Marking a position $\mapsto c$	$\partial_w \mid C$	$_{bn+d} = (n+1)A_{bn+d}$					
	$\overline{\mathcal{C}=\operatorname{Pos}(\mathcal{A})}$		$C(z,w) = \frac{w}{h} \frac{\partial}{\partial w} (A(z,w))$					
	Adding a key $\mapsto \int$	C	$C_{bn+d} = A_{bn+d-1}$					
	$\overline{\mathcal{C}=\operatorname{Add}(\mathcal{A})}$		$ig \ C(z,w) = \int_0^z A(u,w) du$					
	$Bucketing\mapstoexp$	C	$m_{m,n}=\delta(m,1)$					
	$\overline{\mathcal{C}} = \operatorname{Bucket}(\mathcal{Z})$		$(z,w) = w^b \exp(z)$					
	Marking a key $\mapsto \partial_z$	C	$m,n = nA_{m,n}$					
	$\overline{\mathcal{C} = \operatorname{Mark}(\mathcal{A})}$		$(z,w)=zrac{\partial}{\partial z}A(z,w)$					
	$n \mapsto [n] =$	1+q	$\overline{q^2+\ldots+q^{n-1}}=rac{1-q^n}{1-q}.$					
\sum	$\int (n+1)f_n z^n \mapsto \sum [n]$	$(+ 1]f_{1}$	$_{n}z^{n}$.					
$\frac{u}{d}$	$rac{w}{b}rac{\partial}{\partial w}A(z,w) \hspace{2mm}\mapsto \hspace{2mm} H[A(z,w)]$	z,w)]	$t=rac{A(z,w)-A(z,wq^{rac{1}{b}})}{1-q}$					
	$z rac{\partial}{\partial z} A(z,w) \hspace{0.2cm}\mapsto \hspace{0.2cm} \hat{H}[A(z,w)]$	z,w)]	$x=rac{A(z,w)-A(qz,w)}{1-q}.$					

Two models to analyze the problem.

• Exact filling model.

 A fixed number of keys n, are distributed among m locations, and all mⁿ possible arrangements are equally likely to occur.

Poisson model.

 Each location receives a number of keys that is Poisson distributed with parameter bα, and is **independent** of the number of keys going elsewhere. This implies that the total number of keys, N, is itself a Poisson distributed random variable with parameter bαm:

$$Pr[N=n] = rac{e^{-blpha m}(blpha m)^n}{n!}.$$

Parking, random graphs and random trees.



[Spencer 1997]



- BFS traversal of a random graph γ with vertices {0, 1, ..., n}.
- Induces a queue $(H_i(au))_{1\leq i\leq n}$ and a spanning tree au .
- BFS induces a parking sequence $(X_i(au))_{1\leq i\leq n}$.
- Ex: $(X_i(\tau)) = \{\{6, 8\}, \{2, 3\}, \phi, \{7\}, \{1, 4\}, \{5\}, \{9\}, \phi, \phi\}.$
- $x_i(\tau) = |X_i(\tau)|$. <u>Ex:</u> $(x_i(\tau)) = \{2, 2, 0, 1, 2, 1, 1, 0, 0\}$.
- $y_i(\tau) = x_1(\tau) + x_2(\tau) + \ldots + x_i(\tau) i + 1$, size of queue $(H(\tau))$ before step *i*. <u>Ex:</u> $(y_i(\tau)) = \{2, 3, 2, 2, 3, 3, 3, 2, 1\}$.
- $y_1(\tau) + \ldots y_n(\tau) n$ is the total displacement. Ex: 12.
- BFS induces a random walk excursion. <u>Ex:</u> b.

- Let X_i be the number of elements that have hash address i.
- Let H_i be the total number of elements that try bucket i.
- let Q_i be the *overflow* from bucket i.
- We thus have the equations

$$H_i=X_i+Q_{i-1}, \qquad \qquad Q_i=(H_i-b)_+.$$

- The number of elements stored in bucket i is $Y_i = \min(H_i, b)$.
- The bucket is full if and only if $H_i \ge b$.

Finite and infinite hash tables.

Lemma

Let X_i , $i \in \mathfrak{T}$, be given non-negative integers. If $\mathfrak{T} = \{1, \ldots, m\}$ or \mathbb{N} , then the equations (??), for all $i \in \mathfrak{T}$, have a unique solution given by, considering $j \ge 0$,

$$H_i = \max_{j < i} \sum_{k=j+1}^i (X_k - b) + b, \qquad Q_i = \max_{j \leq i} \sum_{k=j+1}^i (X_k - b)$$

• Properly defined, the result can be extended to infinite tables.

Convergence to an infinite hash table.

- We are interested in hashing on Z_m with n elements having independent uniformly random hash addresses.
- X₁,..., X_m have a multinomial distribution with parameters n and (1/m,..., 1/m). (We denote these X_i by X_{m,n;i}.)
- We denote the profile of this hash table by H_{m,n;i}, where as above i ∈ Z_m but we also can allow i ∈ Z in the obvious way.
- We consider a limit with $m,n
 ightarrow\infty$ and $n/bm
 ightarrowlpha\in(0,1).$
- The appropriate limit object is an infinite hash table on \mathbb{Z} with $X_i = X_{\alpha;i}$ that are independent and identically distributed (i.i.d.) with the Poisson distribution $X_i \sim Po(\alpha b)$.
- We denote the profile of this hash table by $H_{\alpha;i}$.

Lemma

Let $m, n \to \infty$ with $n/bm \to \alpha$ for some α with $0 < \alpha < 1$. Then $(H_{m,n;i})_{i=-\infty}^{\infty} \stackrel{d}{\longrightarrow} (H_{\alpha;i})_{i=-\infty}^{\infty}$.

Poisson Transform.

• Results in one model can be transfered into the other model by the Poisson Transform:

$$P_m[f_{m,n};blpha] = \sum_{n\geq 0} Pr[N=n]f_{m,n} = e^{-blpha m} \sum_{n\geq 0} rac{(blpha m)^n}{n!} f_{m,n}.$$



Inversion Theorem: [Gonnet and Munro 1984]



$$\text{If } \mathbf{P}_m[f_{m,n};b\alpha] = \sum_{k \geq 0} a_{m,k}(bm\alpha)^k \text{ then } f_{m,n} = \sum_{k \geq 0} a_{m,k} \frac{n^k}{(bm)^k}.$$

• The Poisson model is an approximation of the exact filling model when $n, m \to \infty$ with $n/m = b\alpha$ with $0 \le \alpha < 1$.

Diagonal Poisson Transform.





- Let a hash table of size m, with n + 1 keys, and let P be a property for ● (chosen uniformily at random).
- Let $f_{m,n}$ be the result of applying a linear operator (e.g. an expected value) to the probability generating function of P.

$$f_{m,n} = \sum_{i \ge 0} \Pr[\bullet \in \text{cluster of size } i+1] \ f_{i+2,i}$$
$$= \sum_{i \ge 0} \binom{n}{i} \frac{(m-i-2)^{n-i-1}(m-n-2)(i+2)^i}{m^n} f_{i+2,i}.$$

• Then $\mathcal{P}_m[f_{m,n};lpha]=\mathcal{D}_2[f_{n+2,n};lpha]$ with

$${\mathcal D}_c[f_n;lpha] = (1-lpha)\sum_{n\geq 0} e^{-(n+c)lpha} rac{((n+c)lpha)^n}{n!} f_n.$$

• As a consequence,

$$egin{aligned} P(z,w,q) &:= \sum\limits_{m,n \geq 0} p_{m,n}(q) \; w^{bm} rac{z^n}{n!} = rac{\hat{N}_0(z,w,q)}{1-N_0(z,w)}, \qquad ext{with} \ \hat{N}_0(z,w,q) &:= \sum\limits_{d=0}^{b-1} w^{b-d} \; \sum\limits_{i \geq 0} F_{bi+d} \; rac{(zw)^{bi+d}}{(bi+d)!} \; p_{bi+d}(q), \end{aligned}$$

which could be directly derived with the sequence construction.

•
$$\hat{N}_0(z,w,1) = N_0(z,w), \ P(z,w,1) = \Lambda_0(z,w) - 1.$$

• Moreover, by singularity analysis of $\Lambda_0(lpha,y^{1/b}e^{-lpha},q)$ at y= 1,

$$\lim_{m\to\infty}\mathbf{P}_m[Q_{m,n,0}/m^n;b\alpha]=T_0(b\alpha)=\frac{b(1-\alpha)}{\prod_{j=1}^{b-1}\left(1-\frac{T(\omega^j\alpha e^{-\alpha})}{\alpha}\right)}.$$

• As a consequence,

$$\lim_{m\to\infty}\mathbf{P}_m[p_{m,n}(q)/m^n;b\alpha] \hspace{2mm} = \hspace{2mm} T_0(b\alpha)\tilde{N}_0(b\alpha,e^{-\alpha},q).$$

Generating Functions and the Poisson Transform.

 Let P_{m,n}(q) be the generating function of a cumulated value of a RV χ in a hash table of size m with n elements. Let

$$P(z,w,q) = \sum_{m \geq 0} w^{bm} \sum_{n \geq 0} rac{P_{m,n}(q)}{m^n} rac{(bmz)^n}{n!}.$$

• Then, for a fixed 0 < lpha < 1,

$$egin{array}{rcl} P(lpha,y^{1/b}e^{-lpha},q) &=& \displaystyle\sum_{m\geq 0}y^m\left(e^{-bmlpha}\sum_{n\geq 0}rac{P_{m,n}(q)}{m^n}rac{(bmlpha)^n}{n!}
ight)\ &=& \displaystyle\sum_{m\geq 0}y^mP_m\left[rac{P_{m,n}(q)}{m^n};blpha
ight]. \end{array}$$

• Results for the probability generating function in the Poisson Model $(n, m \to \infty, 0 \le n/bm = \alpha < 1)$ can be found by singularity analysis from $P(\alpha, y^{1/b}e^{-\alpha}, q)$. In our problems, the dominant singularity is at y = 1.

The overflow (parking problem).

• Let $N_{m,n,k}$ be the number of tables of length m with n elements and overflow k and

$$\Omega(z,w,q):=\sum_{m\geq 0}\sum_{n\geq 0}\sum_{k\geq 0}N_{m,n,k}w^{bm}rac{z^n}{n!}q^k.$$



Proof.

$$egin{aligned} & extsf{Sketch]} \ \Omega(z,w,q) = 1 + \Omega(z,w,q) rac{w^b e^{zq}}{q^b} + \sum_{s=0}^{b-1} (1-q^{s-b}) \mathcal{O}_s(z,w), \end{aligned}$$

with

$$\mathcal{D}_s(z,w)=rac{F_s(zw)w^{b-s}}{1-N_0(z,w)}.$$

The overflow (parking problem).

• Let $Q_{m,n}$ denote the overflow in a random hash table with m buckets and n keys.

Corollary

$$\mathbb{E}\,Q_{m,n}=m^{-n}\sum_{j=0}^n\sum_{k=1}^{\lfloor j/b
floor}inom{n}{j}(j-kb)k^{j-1}(m-k)^{n-j}.$$

The parking problem.

Theorem

Let $0 < \alpha < 1$. The probability generating functions $\psi_H(z)$ and $\psi_Q(z)$ extend to meromorphic functions given by

$$\psi_{H}(z) = rac{b(1-lpha)(z-1)}{z^{b}e^{lpha b(1-z)}-1}rac{\prod_{\ell=1}^{b-1}(z-T(\omega^{\ell}lpha e^{-lpha})/lpha)}{\prod_{\ell=1}^{b-1}(1-T(\omega^{\ell}lpha e^{-lpha})/lpha)},
onumber \ \psi_{Q}(z) = rac{b(1-lpha)(z-1)}{z^{b}-e^{lpha b(z-1)}}rac{\prod_{\ell=1}^{b-1}(z-T(\omega^{\ell}lpha e^{-lpha})/lpha)}{\prod_{\ell=1}^{b-1}(1-T(\omega^{\ell}lpha e^{-lpha})/lpha)}.$$

Corollary

In the infinite Poisson model, the probability of no overflow from a given bucket is

$$\Pr(Q_i=0)=e^{b\alpha}T_{b-1}(b\alpha).$$

This is the asymptotic probability of success in the parking problem, as $m, n \to \infty$ with $n/m \to \alpha$,

The parking problem.

Corollary For k = 0, ..., b - 1. $\Pr(Y_i=k)=\Pr(H_i=k)=-b(1-lpha)rac{[z^k]\prod_{\ell=0}^{b-1}(z-T(\omega^\elllpha e^{-lpha})/lpha)}{\prod_{\ell=1}^{b-1}(1-T(\omega^\elllpha e^{-lpha})/lpha)}.$ Furthermore, the probability that a bucket is not full is given by $\Pr(Y_i < b) = \Pr(H_i < b) = T_0(b\alpha) = rac{b(1-lpha)}{\prod_{\ell=1}^{b-1} (1-T(\omega^\ell lpha e^{-lpha})/lpha)}$ and thus $\Pr(Y_i = b) = \Pr(H_i > b) = 1 - T_0(b\alpha).$

Robin Hood: an example (b=2).

Keys inserted:

• 36, 77, 24, 79, 56, 69, 49, 18, 38, 97, 78, 10, 58. Hash function

•
$$h(x) = x \mod 10$$
.

а	49	79			24		36	77	18	58
	69	10					56	97	38	78
	0	1	2	3	4	5	6	7	8	9

• What happens when 29 is inserted?

a	29	69	10		24		36	77	18	58
	49	79					56	97	38	78
	0	1	2	3	4	5	6	7	8	9

Properties of Robin Hood Hashing.



- At least one record is in its home bucket.
- The keys are stored in nondecreasing order by hash value, starting at some location k and wrapping around. In our example, k = 5 (the first slot of the third bucket).
- If a fixed rule is used to break ties among the candidates to probe their next probe bucket (eg: by sorting these keys in increasing order), then the resulting table is independent of the order in which the records were inserted. Then, we may insert the elements in any order, and study the behavior of the last one inserted!.

Robin Hood displacement.

- W.I.o.g. we search for a record that hashes to bucket 0.
- We have to probe buckets occupied by the elements that would have gone to the overflow area.
- Then consider collisions with all the elements that hash to 0.
- Let D^{RH} be the displacement of a given element x.
- Let C^{RH} be the number of elements that win over x in the competition for slots in the buckets. Then $D^{\mathsf{RH}} = \lfloor C^{\mathsf{RH}}/b \rfloor$.
- The specification is

$$C^{\mathsf{RH}} = Q_{-1} + V = \mathsf{Overflow} * \mathsf{Mark}(\mathsf{Bucket})$$

of the number $Q_{-1} = Q_{m-1}$ of keys that overflow into 0 and the number V of keys that hash to 0 that win over x.

• The number Q_{m-1} of keys that overflow does not change when the keys that hash to 0 are removed. This is thus independent of V.

Robin Hood displacement.

• We consider the displacement D^{RH} of a marked key •.

$$RH(z,w,q):=\sum_{m\geq 0}\sum_{n\geq 0}\sum_{k\geq 0}CRH_{m,n,k}w^{bm}rac{z^n}{n!}q^k$$

where $CRH_{m,n,k}$ is the number of hash tables of length m with n keys (one of them marked as \bullet) such that \bullet hashes to the first bucket and the displacement D^{RH} of \bullet equals k.

Theorem

$$RH(bz, w, q) = rac{1}{b} \sum_{d=0}^{b-1} C\left(bz, w, \omega^d q^{1/b}
ight) \sum_{p=0}^{b-1} \left(\omega^d q^{1/b}
ight)^{-p},$$

with

$$C(bz,w,q) = rac{w^b(e^{bz}-e^{bzq})}{(1-q)(q^b-w^be^{bzq})}rac{\prod_{j=0}^{b-1}\left(q-rac{T(w^jzw)}{z}
ight)}{\prod_{j=0}^{b-1}\left(1-rac{T(w^jzw)}{z}
ight)}.$$

Robin Hood displacement.

• Note that the expectation of the displacement (but not the variance) is the same for any insertion heuristic. We let $D_{m,n}$ denote the displacement of a random element in a hash table with m buckets and n keys.

Lemma

For linear probing with the Robin Hood, FCFS or LCFS (or any other) heuristic,

$$\mathbb{E}\, D_{m,n} = rac{m}{n}\,\mathbb{E}\, Q_{m,n}.$$

Proof.

For any hash table, and any linear probing insertion policy, the sum of the n displacements of the keys equals the sum of the m overflows Q_i . Take the expectation.

Robing Hood displacement.

۲

Theorem

Let $0 < \alpha < 1$. In the infinite Poisson model, the variable V_{α} , the number of keys that win over the new key C_{α}^{RH} and its Robin Hood displacement D_{α}^{RH} have the probability generating functions

$$\psi_V(q) = rac{1-e^{blpha(q-1)}}{blpha(1-q)}$$

$$egin{aligned} \psi_C(q) &= \psi_Q(q)\psi_V(q) = rac{1-lpha}{lpha} rac{1-e^{blpha}(q-1)}{e^{blpha}(q-1)-q^b} rac{\prod_{\ell=1}^{i-1}(q-\zeta_\ell)}{\prod_{\ell=1}^{b-1}(1-\zeta_\ell)} \ \psi_{\mathsf{RH}}(q) &= rac{1}{b}\sum_{j=0}^{b-1}\psi_C(\omega^j q^{1/b})rac{1-q^{-1}}{1-\omega^{-j}q^{-1/b}}. \end{aligned}$$

As $m, n \to \infty$ with $n/bm \to \alpha$, $V_{m,n} \stackrel{d}{\longrightarrow} V_{\alpha}$, $C_{m,n}^{RH} \stackrel{d}{\longrightarrow} C_{\alpha}^{RH}$ and $D_{m,n}^{RH} \stackrel{d}{\longrightarrow} D_{\alpha}^{RH}$, with convergence of all moments; furthermore, for some $\delta > 0$, the corresponding probability generating functions converge, uniformly for $|q| \leq 1 + \delta$.

Robing Hood displacement.

Corollary

As $m,n
ightarrow \infty$ with $n/bm
ightarrow lpha \in (0,1)$,

$$\mathbb{E} C_{m,n}^{\mathsf{RH}} o \mathbb{E} C_{lpha}^{\mathsf{RH}} = rac{1}{2(1-lpha)} - rac{b}{2} + \sum_{\ell=1}^{b-1} rac{1}{1-\zeta_{\ell}},$$

 $\mathbb{E} D_{m,n}^{\mathsf{RH}} o \mathbb{E} D_{lpha}^{\mathsf{RH}} = rac{1}{2blpha} \left(rac{1}{1-lpha} - b - blpha
ight) + rac{1}{blpha} \sum_{\ell=1}^{b-1} rac{1}{1-\zeta_{\ell}}.$

with

$$\zeta_{\ell} := T(\omega^{\ell} \alpha e^{-\alpha})/\alpha.$$

• Let F_{bi+d} be the number of ways to construct an almost full table of length i+1 and size bi+d with $0 \le d \le b-1$. Then,

$$F_d(u):=\sum_{i\geq 0}F_{bi+d}rac{u^{bi+d}}{(bi+d)!}, \ \ \ N_0(z,w):=\sum_{d=0}^{b-1}w^{b-d}F_d(zw).$$

- In an almost full table the length of the block is marked by w in N₀(bz, w).
- The generating function B(z, w, q) for the block length is

$$B(bz,w,q) = \Lambda_0(bz,w) N_0(bz,wq^{1/b}) = rac{1-\prod_{j=0}^{b-1} \left(1-rac{T(\omega^j z w q^{1/b})}{z}
ight)}{\prod_{j=0}^{b-1} \left(1-rac{T(\omega^j z w)}{z}
ight)}$$

Let $B_{m,n}$ be the length of a random block, chosen uniformly among all blocks in all hash tables with m buckets and n keys. This is the same as the length of the last block in a uniformly random hash table such that the rightmost bucket is not full. Recall that we denote the number of such hash tables by $Q_{m,n,0}$.

Corollary If $0 \le n < bm$, then $\mathbb{E} B_{m,n} = rac{m^n}{Q_{m,n,0}}.$

Proof.

The sum of the block lengths in any hash table is m, and thus the sum of the lengths of all blocks in all tables is $m \cdot m^n$, while the number of blocks ending with a given bucket is $Q_{m,n,0}$ and thus the total number of blocks is $m \cdot Q_{m,n,0}$.

• Let B be the length of the first block, i.e.,

 $B := \min\{i \ge 1 : Y_i < b\} = \min\{i \ge 1 : H_i < b\}.$

• Hence, B is the first positive index i such that the number of elements $S_i = X_1 + \cdots + X_i$ hashed to the i first buckets is less than the capacity bi of these buckets, i.e.,

$$B = \min\{i \geq 1: S_i < bi\}.$$

Theorem

The probability generating function $\psi_B(z):=\mathbb{E}\,z^B$ of B is given by

$$\psi_B(z) = 1 - \prod_{\ell=0}^{b-1} \left(1 - T(\omega^\ell lpha e^{-lpha} z^{1/b})/lpha
ight),$$

for $|q| \leq R$ for some R > 1.

Corollary

The random block length $B = B_{\alpha}$ defined above has expectation

$$\mathbb{E}\,B_{lpha}=rac{1}{T_0(blpha)}$$

and variance,

$$\mathbb{V}[B]_{oldsymbol{lpha}}=rac{1}{b(1-lpha)^2T_0(blpha)}-rac{2}{bT_0(blpha)}\sum_{\ell=1}^{b-1}rac{\zeta_\ell}{(1-\zeta_\ell)(1-lpha\zeta_\ell)}-rac{1}{T_0(blpha)^2}$$

• The length of the block \hat{B}_i containing a given bucket i has a different, size-biased distribution.

Theorem

In the infinite Poisson model, $\hat{B}=\hat{B}_{\alpha}$ has the size-biased distribution

$$\Pr(\hat{B}_lpha=k)=rac{k\Pr(B_lpha=k)}{\mathbb{E}\,B_lpha}=T_0(blpha)k\Pr(B_lpha=k)$$

and thus the probability generating function

$$\psi_{\hat{B}}(q)=T_0(blpha)q\psi_B'(q)=T_0(blpha)q\sum_{\ell=0}^{b-1}\zeta_\ell'(q)\prod_{j
eq\ell}(1-\zeta_j(q)).$$

As $m, n \to \infty$ with $n/bm \to \alpha$, $\hat{B}_{m,n} \stackrel{d}{\longrightarrow} \hat{B}_{\alpha}$ with convergence of all moments; furthermore, for some $\delta > 0$, the probability generating function converges to $\psi_{\hat{B}}(q)$, uniformly for $|q| \leq 1 + \delta$.

Corollary

As $m,n
ightarrow\infty$ with $n/bm
ightarrowlpha\in(0,1)$,

$$\mathbb{E}\,\hat{B}_{m,n} o\mathbb{E}\,\hat{B}_{lpha}=rac{1}{b(1-lpha)^2}-rac{2}{b}\sum_{\ell=1}^{b-1}rac{\zeta_\ell}{(1-\zeta_\ell)(1-lpha\zeta_\ell)}.$$

Theorem

In the exact model, $\hat{B}_{m,n}$ has the size-biased distribution

$$\Pr(\hat{B}_{m,n}=k)=rac{k\Pr(B_{m,n}=k)}{\mathbb{E}\,B_{m,n}}=rac{Q_{m,n,0}}{m^n}k\Pr(B_{m,n}=k).$$

Theorem

For the exact model, as $m, n \to \infty$ with $n/bm \to \alpha$, $B_{m,n} \xrightarrow{d} B_{\alpha}$ with convergence of all moments.

Unsuccessful search.

- In a cluster with n keys, the number of visited buckets in a unsuccessful search, is the same as the one needed to insert the (n + 1)st element.
- Then, the specification Pos(C) (marking the position of this inserted element) leads to

$$egin{aligned} U(bz,w,q) &=& \sum_{m\geq 1} w^{bm} \sum_{n\geq 0} rac{(bmz)^n}{n!} P_{m,n}(q) \ &=& \Lambda_0(bz,w) rac{N_0(bz,w)-N_0(bz,wq^{1/b})}{1-q} \ &=& rac{\prod_{j=0}^{b-1} \left(1-rac{T(\omega^j z wq^{1/b})}{z}
ight) - \prod_{j=0}^{b-1} \left(1-rac{T(\omega^j z w)}{z}
ight)}{(1-q)\prod_{j=0}^{b-1} \left(1-rac{T(\omega^j z w)}{z}
ight)}, \end{aligned}$$

where $P_{m,n}(q)$ is the probability generating function for the displacement of the (n + 1)st inserted element.

FCFS displacement.

- We consider the displacement of a marked key •, which we by symmetry may assume hashes to the first bucket.
- Thus, let

$$FCFS(z,w,q):=\sum_{m\geq 1}\sum_{n\geq 1}\sum_{k\geq 0}FCFS_{m,n,k}w^{bm}rac{z^n}{n!}q^k,$$

where $FCFS_{m,n,k}$ is the number of hash tables of length m with n keys (one of them marked as \bullet) such that \bullet hashes to the first bucket and the displacement D^{FC} of \bullet equals k.

For a given m and n with 1 ≤ n ≤ bm, there are nmⁿ⁻¹ such tables (n choices to select • and mⁿ⁻¹ choices to place the other n − 1 elements).

FCFS displacement.

• Thus, if $d_{m,n}(q)$ is the probability generating function for the displacement of a random key

$$egin{aligned} FCFS(z,w,q) &= \sum_{m\geq 1}\sum_{n=1}^{bm} nm^{n-1}d_{m,n}(q)w^{bm}rac{z^n}{n!} \ &= z\sum_{m\geq 1}\sum_{n=0}^{bm-1}d_{m,n+1}(q)w^{bm}rac{m^nz^n}{n!} \end{aligned}$$

Theorem

$$FCFS(bz, w, q) = b \int_0^z U(bt, we^{z-t}, q) dt$$

up to terms $z^n w^m q^k$ with n > bm.

FCFS displacement.

Proof.

The probability generating function for the displacement of a random key when having n keys in the table is

$$d_{m,n}(q) = rac{1}{n}\sum_{i=0}^{n-1} u_{m,i}(q),$$

Then, for all $m\geq 1$ and $n\geq$,

$$u_{m,n}(q)=(n+1)d_{m,n+1}(q)-nd_{m,n}(q),$$

and so,

$$rac{\partial}{\partial z}FCFS(bz,w,q)-wrac{\partial}{\partial w}FCFS(bz,w,q)=bU(bz,w,q).$$

This differential equation together with the boundary condition F(0, w, q) = 0 leads to the solution.

Specification for FCFS displacement (b=1).

• Let FC(z, w, q) be the generating function for the cost of a successful search in an almost full table when n + 1 elements are inserted and one element • is marked.

• $AF(z) = rac{T(z)}{z}$ where T(z) is the tree function.



• Then, $\langle FC \rangle = \langle Add(AF) \rangle^* \langle AF \rangle + \langle Pos(FC) \rangle^* \langle AF \rangle + \langle Pos(AF) \rangle^* \langle FC \rangle$ leads to

$$egin{aligned} \partial_z(FC) &= H[AF] * AF + rac{\partial}{\partial z}(zFC) * AF + rac{\partial}{\partial z}(zAF) * FC. \ FCFS(z,w,q) &= rac{((1-T(zwq))^2-(1-T(zw))^2)T(zw)}{2z(1-q)(1-T(zw))}. \end{aligned}$$

Unsuccessful search and FCFS displacement.

Theorem

The probability generating function $\psi_U(z) := \mathbb{E} \, z^{U_i}$ of U_i is given by

$$\psi_U(z) = rac{T_0(blpha)}{1-z} \prod_{\ell=0}^{b-1} \left(1-Tig(\omega^\ell lpha e^{-lpha} z^{1/b}ig)/lphaig).$$

Theorem

The probability generating function $\psi_{FC}(z) := \mathbb{E} z^{D_i^{FC}}$ of D_i^{FC} is given by

$$egin{aligned} \psi_{\mathsf{FC}}(z;lpha) &= rac{1}{lpha} \int_0^lpha \psi_U(z;eta) \, \mathrm{d}eta &= rac{1}{lpha} \int_0^lpha rac{ au(eta)}{1-z} \prod_{\ell=0}^{b-1} (1-\zeta_\ell(z;eta)) \, \mathrm{d}eta \ &= rac{1}{lpha} \int_0^lpha rac{b(1-eta) \prod_{\ell=0}^{b-1} (1-\zeta_\ell(z;eta))}{(1-z) \prod_{\ell=1}^{b-1} (1-\zeta_\ell(1;eta))} \, \mathrm{d}eta. \end{aligned}$$

Some final considerations.

- Problem with a very rich history.
- Paradigm of a problem that nicely integrates analytical, combinatorial and probabilistic approaches.
- This integration (together with the use of symbolic methods!) has allowed the understanding of deep relations with other important problem.
- A unified analysis of several important random variables related with linear probing: symbolic methods + random walks (linked by the Poisson Transform).

Some ongoing work.

- Total displacement with buckets. Relation with other problems as for b=1?
- Number of movements in deletion algorithm.