

ANALYSIS OF LINEAR
PROBING HASHING WITH
BUCKETS
(SURVEY)

ALFREDO VIOLA

(viola@ting.edu.uy)

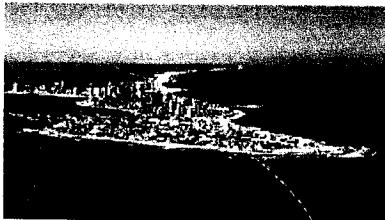
- Universidad de la Repùblica, Montevideo, URUGUAY
- LIPN - Université de Paris-Nord, Villejuif, FRANCE



Announcement and Call for Papers

Information Theory Workshop (ITW'2006)

March 13–17, 2006, Punta del Este, Uruguay



General Co-chairs

Gadiel Seroussi
Alfredo Viola

Program Committee

Ron Roth, co-chair
Marcelo Weinberger, co-chair
Mario Blaum
Michelle Effros
Ioannis Kontoyiannis
Hugo Krawczyk
Gábor Lugcs
Muriel Médard
Neri Merhav
Alon Orlitsky
Daniel Panario
Serap Savari
Amin Shokrollahi
Wojciech Szpankowski
Ruediger Urbanke
Alexander Vardy
Sergio Verdú

Organizing Committee

Pablo Belzarena
Gustavo Brown
Álvaro Martín
Ignacio Ramírez

Important Dates

Submission: Oct. 31, 2005
Notification: Dec. 23, 2005
Final version: Jan. 15, 2006

The 2006 IEEE Information Theory Workshop (ITW'2006), organized by the Sociedad de Ingeniería—Universidad de la República, Uruguay and the IEEE Information Theory Society, will take place from the evening of Monday March 13, 2006 to Friday, March 17, 2006, at the Conrad Resort and Casino in Punta del Este, Uruguay.

The sessions of the workshop will cover the following topics:

- Algebraic and combinatorial coding theory
- Algorithms in finite fields
- Analysis of algorithms in information theory
- Application of coding theory in computer science
- Coding techniques for storage
- Communication complexity
- Cryptography and data security
- Data compression algorithms and source coding
- Data networks
- Detection
- Information theory and statistics
- LDPC codes, turbo codes, and iterative decoding
- Multi-user information theory
- Pattern recognition and learning
- Shannon theory
- Universal prediction and on-line algorithms

Sessions will consist mainly of invited talks, but will also include slots for contributed papers. Submission of papers on the above topics is hereby solicited. The deadline for submission is October 31, 2005. Authors will be notified of acceptance decisions by December 23, 2005. The final versions of all contributions, to be published in the workshop proceedings CD, will be due by January 15, 2006.

Further information, such as submission guidelines, contacts, and local information, will be available at the workshop web site,

<http://www.fing.edu.uy/itw06>.

Location: Punta del Este is one of the premier resort cities in South America. It offers cultural attractions, sophistication, excellent restaurants, natural beauty, and fabulous beaches (the workshop takes place at the end of summer). It is about 1.5 hrs of scenic drive from the capital Montevideo, or a 30 min flight from Buenos Aires, Argentina.

Announcement and Call for Papers
Latin American Theoretical INformatics (LATIN'2006)
March 20–24, 2006, Valdivia, Chile

The 7th Latin American Theoretical INformatics Symposium will take place in March, from Monday 20 through Friday 24, 2006, in Valdivia, Chile. Submissions in all areas of theoretical computer science are welcome, including (but not limited to):

- Algorithms
- Automata theory and formal languages
- Coding theory and data compression
- Combinatorics and graph theory
- Complexity theory
- Computational algebra
- Computational biology
- Computational geometry
- Computational number theory
- Cryptography
- Databases and information retrieval
- Data structures
- Internet and the web
- Logic in computer science
- Machine learning
- Mathematical programming
- Parallel and distributed computing
- Pattern matching
- Quantum computing
- Random structures and algorithms
- Scientific computing

Sessions will consist of keynote addresses and contributed papers. Submission of papers on the above topics is hereby solicited. The final versions of all accepted contributions will be published in the symposium proceedings by Springer-Verlag, in the Lecture Notes in Computer Science Series.

Keynote addresses will be given by:

- Ricardo Baeza-Yates, U. Chile
- Anne Condon, U. British Columbia
- Ferran Hurtado, U. Politècnica de Catalunya
- R. Ravi, Carnegie Mellon U.
- Madhu Sudan, MIT
- Sergio Verdú, Princeton U.
- Avi Wigderson, Institute for Advanced Study

Further information, such as submission guidelines, contacts, and local information, will be available at the conference web site,

<http://www.latin06.org>.

Important Dates

Submission: Sep. 21, 2005
Notification: Nov. 21, 2005
Final version: Dec. 16, 2005

Location: Valdivia is located at the confluence of three rivers, 15 km (10 mi) from the coast and 340 km (215 mi) south of Santiago. It is a university town of 140,000 inhabitants, also called the "City of rivers" and "Southern Pearl." Recognized by its uniqueness due to its architectural and natural beauty, it harmonically combines Spanish and German heritage. The city proudly exhibits historical fortifications, a botanical garden, and reputed beer industry in a gorgeous natural setting.

DISTRIBUTIONAL
ANALYSIS OF
ABBIN' HOD LINEAR PROBING HASHING
WITH BUCKETS

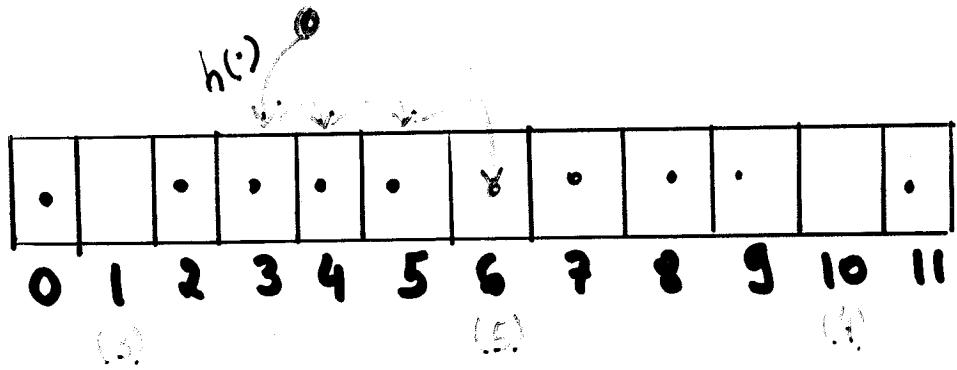
ALFREDO VIOLA

(viola@fing.edu.uy)

- Universidad de la Repùblica, Montevideo, URUGUAY.
- LIPN - Université de Paris-Nord, Villejuif, FRANCE.

MOTIVATION

- Linear Probing is the simplest collision resolution scheme for open addressing.
- Works well when the table is not full.
- Performance deteriorates when load factor increase.
- First proposed by Peterson in 1957.
- The analyses of problems related with linear probing show relations with other problems like tree inversions, tree path lengths, graph connectivity, area under excursions, etc.



FURTHER RESULTS

Individual Displacements:

- (Knuth 1962-1963) First nontrivial algorithm he analyzed.
- (Knuth-Wenm 1966, Rodden-Vieira 1967) First published analyses.
- (Janson, Vieira To appear) Distributional Analysis for LCFS, FCFS, RH
2005-2006
- (Fittell 1925, Johnson-Vieira ?) Longest probe

Construction (cont.):

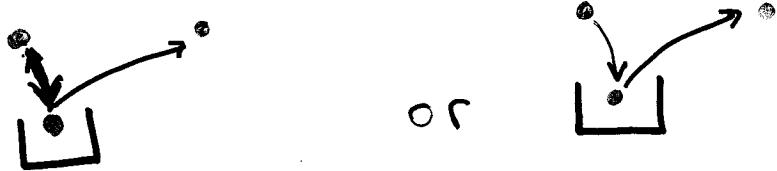
- (Knuth, Flajolet-Poblete-Vieira 1993) Distributional Analysis, and relation with graph connectivity, tree inversions, tree path lengths, area under excursions, etc. ("Airy phenomena")
- (Janson 2001) More limit distributions.
- (Chassaing-Louchard 2002) Phase Transitions.

Buckets with capacity $k \geq 1$. Individual Displacements

- (Blalock-Konheim 1977) Exact PGF for FCFS and asymptotic analysis of the expected value.
- (Poblete-Vieira 1992) Exact expected value for full tables.
- (Poblete-Vieira 1992) Expected value for α -full tables
- (Knuth-Vieira 1963) Distributional Analysis of Robin Hood
- (Vieira 2005) Distributional Analysis of Robin Hood
(NEW)

CRITICAL OBSERVATION WITH IMPERFECT CONSISTENCY

- When two keys collide, either key can get the location. This can be done without looking ahead.



- Several policies can be considered without looking ahead
 - FCFS standard (Poblete - Munro 1989)
 - LCFS (Celis - Larson - Munro 1985)
 - Robin Hood (Celis - Larson - Munro 1987) (Carlsson - Munro - Poblete 1987) - da Linha Probing.
- Neither of these heuristics change the expected value.
- Neither of these heuristics minimizes the variance among all heuristics for Linear Probing that do not look ahead.
- $C_{b,m,n} \rightarrow$ RV for number of probes in a successful search of a hash table with m buckets of capacity b and n elements.
- $C_{b,d} = \lim C_{b,m,n}$

LAW OF LARGE RESULTS FOR A BIN WORD

- $C_{m,n,b} \rightarrow RV$ for cost of successful search of a random element when we insert $m+1$ elements in m buckets of capacity b .

(Knuth 1963, Konheim-Wein 1966)

$$E[C_{m,n}] = \sum_{i=1}^m P_i(m,n) \cdot i$$

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

- $E[C_{m,\alpha m}] = \frac{1}{2}\left(1 + \frac{1}{1-\alpha}\right) - \frac{1}{2(1-\alpha)^3 m} + O\left(\frac{1}{m^2}\right)$

- $E[C_{m,m,1}] = \frac{\sqrt{2\pi m}}{4} + \frac{1}{3} + \frac{1}{48}\sqrt{\frac{2\pi}{m}} + O\left(\frac{1}{m}\right)$

(Catalan - Minimax - Poblete 1977)

- $\text{Var}[C_{m,n,1}] = \frac{1}{2}Q_1(m,n) - \frac{1}{4}Q_0(m,n)^2 - \frac{1}{6}Q_0(m,n) + \frac{n}{6m} - \frac{1}{12}$

- $\text{Var}[C_{m,\alpha m}] = \frac{1}{4(1-\alpha)^2} - \frac{1}{6(1-\alpha)} - \frac{1}{12} + \frac{\alpha}{6} - \frac{1}{6m} - \frac{1+2\alpha}{3(1-\alpha^2)m} + O\left(\frac{1}{m^2}\right)$

- $\text{Var}[C_{m,m,1}] = \frac{4-\pi}{8}m + \frac{1}{9} - \frac{\pi}{48} + \frac{1}{135}\sqrt{\frac{2\pi}{m}} + O\left(\frac{1}{m^2}\right) \rightarrow \text{OPTIMAL!}$

(Janson, Viola To Appeal
2005-2006)

- $C_{\alpha,n}(z) = z \frac{1-\alpha}{\alpha} \frac{e^{2\alpha} - e^{\alpha}}{ze^{\alpha} - e^{2\alpha}}$

- $C_{m,n,1}(z) = \sum_{r=1}^{n+1} \left(\sum_{i=r}^{n+1} \binom{n+1}{i} \frac{(m-n-1+i)}{(n+1)m^i} \sum_{k=0}^{i+r} (-1)^{i+1-r-k} \binom{i+1}{r+k} k^i \right) z^r$

- $P_1\left(\frac{C_{m,m-1,1}}{\sqrt{m}} \leq x\right) \rightarrow \text{Rayleigh}(1/2)$

(Knuth Vol 3)

- $E[C_{m,\alpha m,b}] = 1 + \sum_{k \geq 1} e^{-bkx} \sum_{j \geq 1} j \frac{(bk)^{bk+j-1}}{(bk+j)!}$

(Poblete-Viola 1998)

- $E[C_{m,bm-1,b}] = \frac{1}{b} \sum_{i \geq 2} \binom{bm-1}{i} \frac{(-b)^i}{m^i} \sum_{k=1}^m k^{i-1} \binom{bk-i}{bk-1} + \frac{m-1}{2bm} + 1$

- $E[C_{m,bm-1,b}] = \frac{1}{b} \sum_{i \geq 2} \binom{bm-1}{i} \frac{(-b)^i}{m^i} \sum_{k=1}^m k^{i-1} \binom{bk-i}{bk-1} + \frac{m-1}{2bm} + 1$

Robin Hood - Example

- Insert records with keys $36, 77, 24, 79, 56, 69, 49, 18, 38, 97, 78, 10, 40, 70$ in a table with 10 buckets of size 2 and $h(x) = x \bmod 10$

58

69	10	70	24	36	77	18	78		
79	40			56	97	38	49		
0	1	2	3	4	5	6	7	8	9

- Insert 58.

29

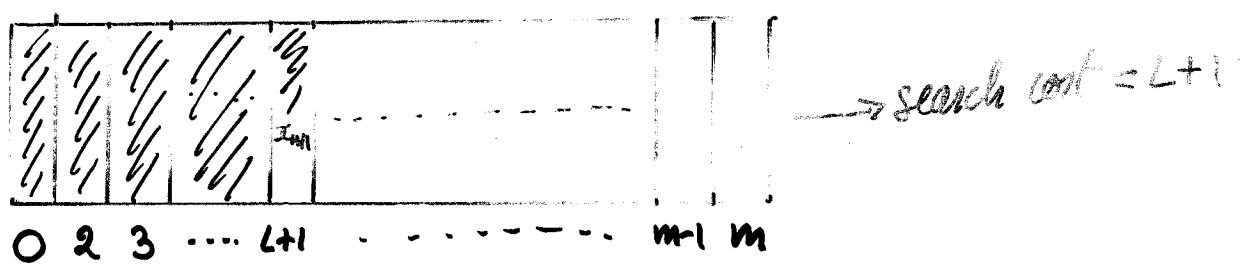
49	79	40	24	36	77	18	58		
69	10	70		56	97	38	78		
0	1	2	3	4	5	6	7	8	9

- Insert 29.

29	69	10	70	24	36	77	18	58	
49	79	40			56	97	38	78	
0	1	2	3	4	5	6	7	8	9

PROPERTIES OF ROBIN HOOD

- At least one record is in its home bucket ($2, 4, 6, 8, 10, 12, 14, 16$)
- If a fixed rule is used to break ties (eg: lexicographical order), then the final table is the same regardless of the order of insertion.
 - ∴ We may search for the last element inserted with $h(x_m) = c_j$
- The keys are stored in non-decreasing order by hash value, starting at some location h and wrapping around (From key 10)



In the first L buckets we have two kinds of records

- Records that overflow
- Records that hash to 0 and won't tie

- We study just the overflow area (Linear Probing Sort) and then the search cost. (Gennet-Nelson 1974, Wilson-Poblete-Thakur 1997, Vidy-Poblete 1997)

POISSON TRANSFORM

- $\mathcal{S}_m [f_{m,n}; b\alpha] = e^{-mb\alpha} \sum_{n \geq 0} \frac{(b\alpha)^n}{n!} f_{m,n} = \sum_{k \geq 0} a_{m,k} (b\alpha)^k$
- $f_{m,n} = \sum_{k \geq 0} a_{m,k} \frac{n^k}{m^k}$

DISTRIBUTION OF BUCKET OCCUPANCY

$Q_{m,n,d} \rightarrow$ # of ways to insert n records in a table with m buckets of size b so that a given bucket has more than d empty slots.

$$\frac{Q_{m,n,b-d+1}}{m^n} - \frac{Q_{m,n,b-d}}{m^n} = p_n \{ \text{given bucket has exactly } d \text{ elements} \}$$

- $Q_{0,n,d} = [n=0]$
- $Q_{m,n,d} = \begin{cases} \sum_{j=0}^n \binom{n}{j} Q_{m-1,j,d} & (0 \leq n < bm+d) \\ 0 & (n \geq bm+d) \end{cases}$

- If $[Q_{m,d}(z)]_n = \sum_{k \geq 0} Q_{m,k,d} \frac{z^k}{k!}$ we have

$$\left\{ \begin{array}{l} Q_{0,d}(z) = 1 \\ Q_{m,d}(z) = \underbrace{[e^z Q_{m-1,d}(z)]}_{(m \geq 1)} \Big|_{bm-d-1} \end{array} \right.$$

Key idea!! \rightarrow New sequence of numbers

If we do NOT have truncation then

$$\left\{ \begin{array}{l} Q_{0,d}(z) = 1 \\ Q_{m,d}(z) = e^z Q_{m-1}(z) \quad (m > 0) \end{array} \right\} \Rightarrow Q_{m,d}(z) = e^{mz}$$

Truncation complicates the solution and it does not seem to be a "nice" closed formula for $Q_{m,d}(z)$.

KEY IDEA! Find a sequence of numbers $T_{k,d,b}$ with $k \geq 0$

and $T_{d,b}(z) = \sum_{k \geq 0} T_{k,d,b} \frac{z^k}{k!}$, INDEPENDENT of m

such that
$$Q_{m,d}(z) = [T_{d,b}(z) e^{mz}]_{bm=d-1}$$

Moreover satisfies

$$\lim_{m \rightarrow \infty} S_m \left[\frac{Q_{m,n,d}}{m^n}, b\alpha \right] = T_{d,b}(b\alpha)$$

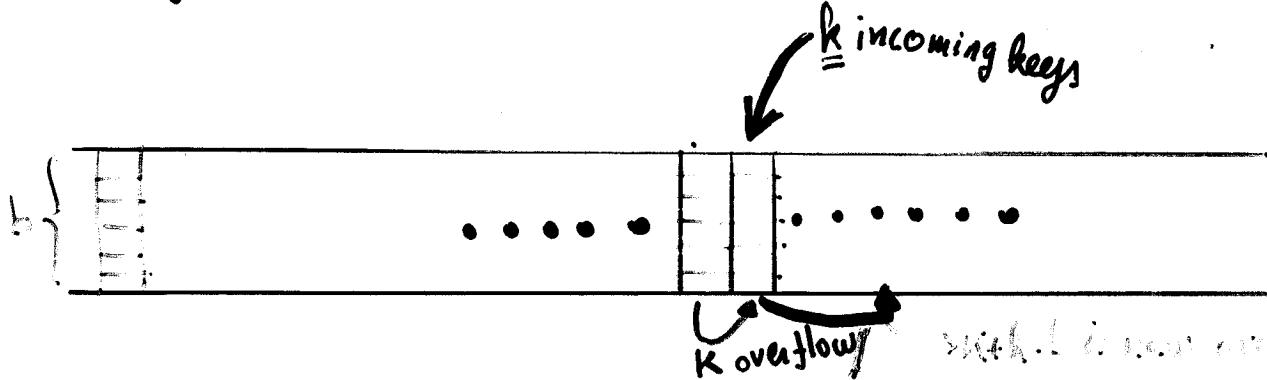
The family $T_d(b\alpha)$ with $0 \leq d \leq b-1$ captures the distribution of bucket occupancy in the Poisson Model!!!

The family of sequences $T_{k,d,b}$

- $T_{k,0,1} \rightarrow 1, -1, 0, 0, 0, 0, \dots$
- $T_{k,0,2} \rightarrow 1, 0, -1, 2, -8, 48, -378, 3672, -42368, \dots$
- $T_{k,1,2} \rightarrow 1, -1, 1, -2, 8, -48, 378, -3672, 42368, \dots$
- $T_{k,0,3} \rightarrow 1, 0, 0, -1, 3, -6, -12, 264, -2016, \dots$
- $T_{k,1,3} \rightarrow 1, 0, -1, 2, -3, -2, 60, -408, 1341, \dots$
- $T_{k,2,3} \rightarrow 1, -1, 1, -1, 0, 8, -48, 144, 675, \dots$
- These sequences do not appear in Sloane's Encyclopedia.
- Some interesting properties
- $\sum_{j=0}^k \binom{k}{j} \left(\lfloor \frac{k+d}{b} \rfloor\right)^{k-j} T_{j,d,b} = [k=0] \Rightarrow \text{implicit definition!}$
- $\sum_{d=0}^{b-1} T_{k,d,b} = \begin{cases} b & (k \geq 0) \\ -1 & (k=1) \\ 0 & (k > 1) \end{cases} \Rightarrow \sum_{d=0}^{b-1} T_{d,b} (b\alpha) = b(1-\alpha)$
- $T_{0,d,b} = 1 \quad T_{k,d,b} = 0 \text{ if } 1 \leq k \leq b-d-1$
- Blaked Konheim (1977) present $T_{0,b}(\alpha)$ and gives some properties.
- Much more work to be done with $T_{k,d,b}!!!$

Poisson Analysis of the overflow area

- We have $m, n \rightarrow \infty$ and $\frac{n}{bm} = \alpha$ ($0 \leq \alpha < 1$)
- $P_k [k \text{ keys hash to a given bucket}] = e^{-b\alpha} \frac{(b\alpha)^k}{k!}$



- $\mathcal{R}(b\alpha, z) \rightarrow \text{PGF for \#elements that overflow from a bucket.}$

First approach: $\mathcal{R}(b\alpha, z) \sim \frac{\mathcal{R}(b\alpha, z) e^{b\alpha(z-1)}}{z^b}$ NOT valid if bucket receives less than b elements

Correction term: $\sum_{s=1}^b (1-\bar{z}^s) P_s(b\alpha)$ \Rightarrow Shaving off elements in bucket.

$$= T_{s-1}(b\alpha) - T_s(b\alpha)$$

After calculations:
$$\boxed{\mathcal{R}(b\alpha, z) = \frac{z-1}{z^b - e^{b\alpha(z-1)}} \sum_{s=0}^{b-1} T_s(b\alpha) \bar{z}^s}$$

When $b=1$ $\mathcal{R}(\alpha, z) = \left(\frac{z-1}{z - e^{\alpha(z-1)}} \right) (1-\alpha)$

→ Eulerian numbers!

(Числа Эйлеровы в интерпретации?)

P.G.F. FOR THE SEARCH COST

- We search for an element pushed by
 - overflow from previous bucket
 - other elements hashed in same bucket
- Let $C(b\alpha, z) = \sum_{k \geq 0} C_k(b\alpha) z^k$ the P.G.F. for the total displacement without considering the fact that we should only count # buckets.

So, the P.G.F. we need is

$$P(b\alpha, z) = z \sum_{k \geq 0} C_k(b\alpha) z^{\lfloor \frac{k}{b} \rfloor} = \sum_{k \geq 0} \left(\sum_{d=0}^{b-1} C_{b\alpha+d}(b\alpha) \right) z^k$$

$$T(b\alpha, z) = \frac{z}{b} \sum_{j=0}^{b-1} C(b\alpha, e^{\frac{2\pi i}{b} j} z^{\frac{1}{b}}) \sum_{p=0}^{b-1} (e^{\frac{2\pi i}{b} p} z^{\frac{1}{b}})^{-p}$$

- If k elements collide with the searched one, the total displacement (expected) originated by separately retrieving all these records is:

$$\frac{1}{k+1} \sum_{r=0}^k z^r = \frac{1}{k+1} \frac{1-z^{k+1}}{1-z}$$

$$\text{with P.G.F } \frac{e^{-b\alpha}}{1-z} \sum_{k \geq 0} \frac{(b\alpha)^k}{(k+1)!} (1-z^{k+1}) = \frac{1-e^{b\alpha(z-1)}}{b\alpha(1-z)}$$

- Then, we have

$$(1 - \frac{1-e^{b\alpha(z-1)}}{b\alpha(1-z)})^{-1} = \frac{e^{b\alpha(z-1)} - 1}{b\alpha(1-z)} \sum_{s=0}^{b-1} T_{b-s}(b\alpha) z^s$$

DISTRIBUTION OF SEARCH COST

Let $T_{b\alpha}$ the R.V. for the cost of searching a random element in a box-full table with buckets of size b using the Robin Hood linear probing algorithm.

$$P_n(b\alpha) = P\{T_{b\alpha} = n+1\} = \sum_{j=0}^{b-1} \frac{T_{b+j-1}(b\alpha)}{b\alpha} e^{-b\alpha} - \sum_{d=0}^{b-1} \left(\frac{(kb\alpha)^{b(n+k)+d-1}}{(b(n+k)+d-1)!} - \frac{(kb\alpha)^{b(n+k+1)+d-1}}{(b(n+k+1)+d-1)!} \right)$$

For fixed m, n the result follows from de poissonization. EXACT!

First moment

$$E[T_{b,m,n+1}] = 1 + \sum_{k=1}^{\lfloor n/b \rfloor} \sum_{i=kb}^n (-1)^{i-kb} \binom{i-1}{kb-1} \frac{(kb)^i}{(i+1)!} \frac{n^i}{(bm)^i}$$

$$E[T_{b\alpha}] = 1 + \sum_{k \geq 1} e^{-b\alpha k} \sum_{m \geq 1} n \frac{(b\alpha k)^{n+bk-1}}{(n+bk)!}$$

$$bE[T_{b,m,bm-1}] = \frac{\sqrt{2\pi bm}}{4} + \frac{1}{3} + \sum_{d=1}^{b-1} \frac{1}{1 - T(e^{\frac{2\pi i}{b} d-1})} + \frac{1}{48} \sqrt{\frac{2\pi}{bm}} + O\left(\frac{1}{bm}\right)$$

With $T(x) = \int_0^x e^{-t^2/2} dt$

HIGHER MOMENTS

$$E\left[T_{ba}^r\right] = r \sum_{s=0}^{b-1} \frac{T_{b+s}(ba)}{ba} \sum_{n \geq 1} \sum_{k \geq 1} e^{-kba} \sum_{d=0}^{b-1} \frac{(bk\alpha)^{b(k+n)+d-1}}{(b(k+n)+d-s)!}$$

$$+ 1 \sum_{s=0}^{b-1} \frac{T_{b+s}(ba)}{ba} \sum_{k \geq 1} e^{-bk\alpha} \sum_{d=0}^{b-1} \frac{(bk\alpha)^{bk+d-1}}{(bk+d-s)!}$$

- For exact m, n we do deoissonization.

Limit distributions

- With $\alpha \neq 1$:
- When $\alpha \rightarrow 1$ main asymptotic contribution of $E\left[T_{ba}^r\right]$ comes from r^{th} derivatives of $C(b\alpha, z^{1/b})$ in $T(b\alpha, z)$.
- As a consequence, proper generalizations of limit distributions for $b=1$ (χ^2 , t , β , Γ) apply:

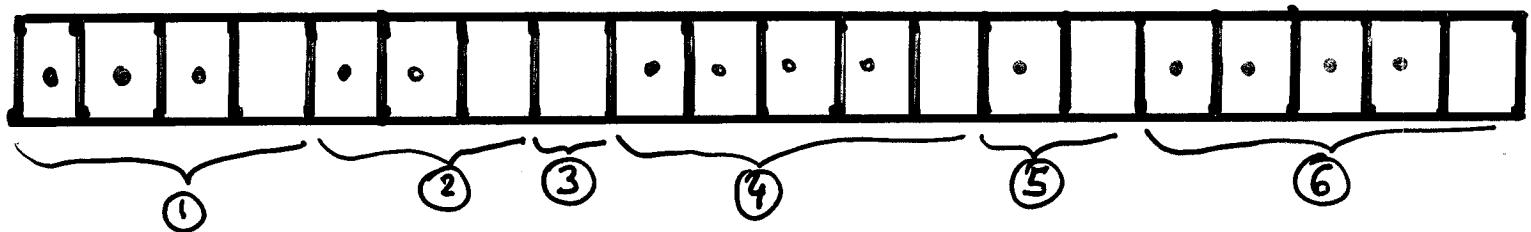
- | |
|---|
| When $\alpha \rightarrow 1$ $b \ln \{ (1-\alpha) T_{ba} \leq x \} \rightarrow$ Exponential ($b/2$) |
| When $m \rightarrow \infty$ $b \ln \left\{ \frac{T_{m,bm-1}}{\sqrt{bm}} \leq x \right\} \rightarrow$ Rayleigh ($1/2$) |

TOTAL CONSTRUCTION COST

Number of filled slots in a cluster $\rightarrow b=1$

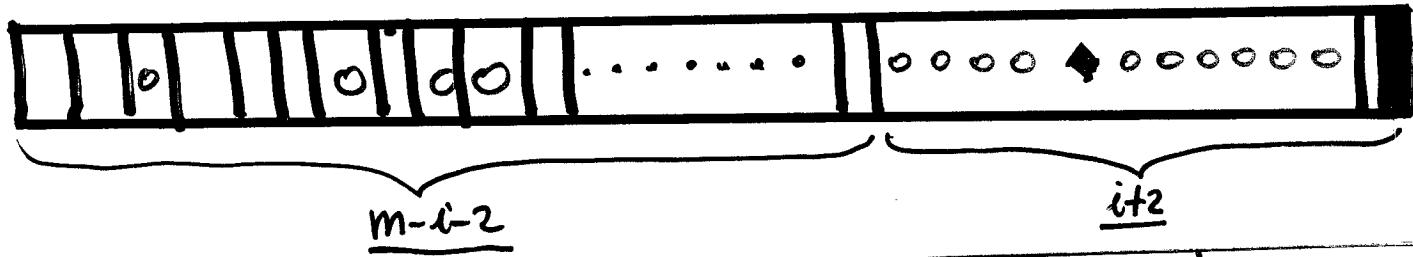
COMBINATORIAL INTERPRETATION

- A linear probing hash table is a sequence of "almost full tables" (clusters).

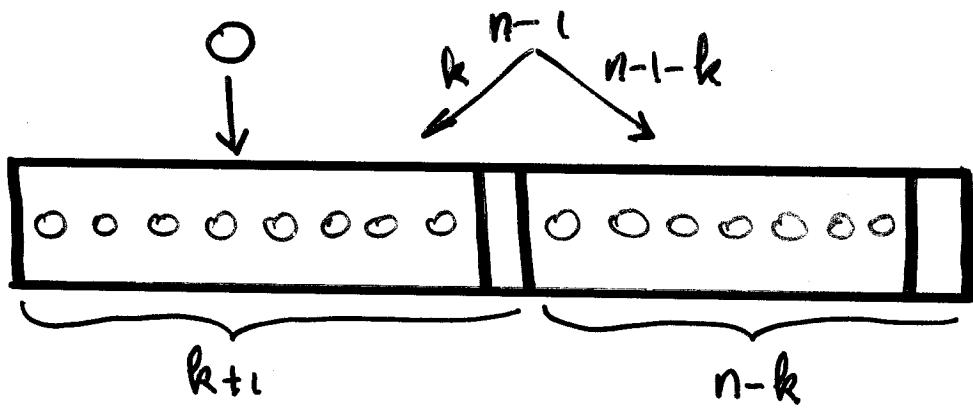


SO, ENOUGH TO STUDY ALMOST FULL TABLES!

- (Hwang, F. & H. S. Kim. 1997.) \rightarrow The Poisson Transform applied to a R.V. in a full table, is the same as the Diagonal Poisson Transform applied to the same RV studied in an almost full table!!! \rightarrow For general b ?



$$m^{\alpha} = (m\alpha)^n, \quad \alpha \sim \beta e^{-\alpha} \frac{\alpha^{-(i+2)\alpha} ((i+2)\alpha)^i}{i!} t_{i+2} i!$$



- $F_{n,k} \rightarrow$ # ways of creating an AFT with n elements and total displacement k

$$\bullet F(z,q) = \sum_{n,k \geq 0} F_{n,k} q^k \frac{z^n}{n!}$$

$$\bullet F_n(q) = \sum_{k=0}^{n-1} \binom{n-1}{k} F_k(q) (1+q+\dots+q^k) F_{n-1-k}(q)$$

$$\frac{\partial}{\partial z} F(zq) = F(zq) \frac{F(z,q) - q F(zq,q)}{1-q}$$

$\rightarrow H.F(z,q)$

- If $d_{n,n-1}$ is the R.V. for the total displacement in a AFT then

$$P_n \left\{ \frac{d_{n,n-1}}{\left(\frac{n}{2}\right)^{3/2}} \leq x \right\} \xrightarrow{(n \rightarrow \infty)} P_1 \left\{ X \leq x \right\}$$

↓ Airy

b $\xrightarrow{bk+b-1}$ $\xrightarrow{bn+d-1}$ $b(n-k)+d$

\circ							
\circ							
\circ							
\circ							
\circ							

- $F_d(q) = 1$

- $F_{bn+d}(q) = \sum_{k=0}^{n-1} \binom{bn+d-1}{bk+d-1} F_{bk+b-1}(q) \frac{1-q^{k+1}}{1-q} F_{b(n-k-1)+d}(q)$
 $+ [d > 0] \frac{1-q^{n+1}}{1-q} F_{bn+d-1}(q)$

- $F_d(z, q) = \sum_{n \geq 0} F_{bn+d}(q) \frac{z^{bn+d}}{(bn+d)!}$

- $H F_d(z, q) = \sum_{n \geq 0} \frac{1-q^{n+1}}{1-q} F_{bn+d}(q) \frac{z^{bn+d}}{(bn+d)!}$

$$F_0(z, q) = e^{\int_0^z H F_d(t, q) dt}$$

$$F_d(z, q) = F_0(z, q) \int_0^z \frac{H F_{d+1}(t, q)}{F_0(t, q)} dt \quad (d \leq d < b)$$

Can solve it!!! → limit distributions? → not even moments!

INDIVIDUAL DISPLACEMENTS: FCFS ($b \geq 1$)

- (Babai, Alon, ...)

- Very interesting ideas and results
- Extremely difficult to read !!
- It does NOT use combinatorial interpretation of linear probing
 - ↳ NEW results at the light of interpreting combinatorially the results in the paper !!

$$\bullet T_{b,d}(z) = \sum_{n \geq 0} F_{bn+d} \frac{z^n}{(bn+d)!} = \frac{1}{z^{d/b}} F_d(z^{1/b}, 1) \quad \text{it does NOT use } q\text{-analog}$$

- Given a sequence $\{\alpha_k^{(2)}, 1 \leq k \leq n\}$ the elementary symmetric functions $\{G_n(x)\}$ of the sequence $\{\alpha_k(z)\}$ are defined as the coefficients of the polynomial $\sum_{k=0}^n G_k(z) x^{n-k} = \prod_{k=1}^n (x + \alpha_k(z))$

$$\bullet \text{Given } \{\alpha_d(z) = T(\rho^d z) : 0 \leq d \leq b-1\} \text{ with } \begin{array}{l} T \rightarrow \text{tree function} \\ r = e^{\frac{2\pi i}{b}} \rightarrow b^{\text{th}} \text{ root of unity} \end{array}$$

$$\text{we have } (bz)^b T_{b,d}((bz)^b) = (-1)^{b-d-1} b^{b-d} G_{b-1}(z)$$

$$z^{b-d} F_d(bz, 1) = (-1)^{b-d-1} G_{b-1}(z)$$

↳ Can it be generalized for q -analogues of F_d ??

$$\Lambda(z, w) = \sum_{m \geq 0} \left(\sum_{n=0}^{bm-1} Q_{m,n,0} \frac{z^n}{n!} \right) w^{bm}$$

After lengthy and complicated derivation

$$\Lambda(z, w) = \frac{N(z, w)}{1 - N(z, w)}$$

$$\text{with } N(z, w) = \sum_{d=0}^{b-1} w^{b-d} F_d(wz, 1)$$

$$= 1 - \prod_{d=0}^{b-1} \left(1 - \frac{b}{2} T\left(r^d \frac{wz}{b}\right) \right)$$

"TRIVIAL" proof if we use combinatorial interpretation!

$$\sum_{n \geq 0} F_n \frac{z^n}{n!} = 1 - \prod_{d=0}^{b-1} \left(1 - \frac{b}{2} T\left(r^d \frac{z}{b}\right) \right)$$

If $w=1$ translates into $\hookrightarrow q\text{-analog}?$

NEW!!!

↓
by combinatorial
interpretation!

$$\Lambda_d(z, w) = \sum_{m \geq 0} \left(\sum_{n=0}^{bm-d-1} Q_{m,n,d} \frac{z^n}{n!} \right) w^{bm}$$

$$\Lambda_d(z, w) = \frac{N_d(z, w)}{1 - N(z, w)}$$

$$N_d(z, w) = N(z, w) - \sum_{j=b-d}^{b-1} \frac{1}{z^{b-j}} [x^j] \left(x^b - bT\left(r^j \frac{zw}{b}\right) \right)$$

Displacement of n^{th} ball

- $G_{m,n,j} = \# \text{ ways of inserting the } n^{\text{th}} \text{ ball with displacement } j$

$$G(z, w, q) = \sum_{m \geq 0} w^{bm} \sum_{n=1}^{b^m} \frac{(mz)^{n-1}}{(n-1)!} \sum_{j=0}^n G_{m,n,j} q^j$$

*Have any hash table
with last bucket not full*

$$= \boxed{H(z, w)}$$

$$\sum_{d=0}^{b-1} w^{b-d} H F_d(w z, 1)$$

insert in last bucket!

"TRIVIAL"
with
combinatorial
interpretation!

Generating Function of $T_0(b\alpha)$

$$\lim_{\substack{m, n \rightarrow \infty \\ \frac{n}{m} = b\alpha}} \frac{Q_{m,n,0}}{m^n}$$

they did not
realize that!:

$$T_0(b\alpha) = \frac{b(1-\alpha)}{\prod_{d=1}^{b-1} \left(1 - \frac{T(r^d \alpha e^{-\alpha})}{\alpha}\right)}$$

for general $T_d(b\alpha)$? \rightarrow can use implicit
definition of $T_{k,d}$?

Expected cost of a random element

- $E[C(b\alpha)] \sim \frac{1}{2b(1-\alpha)} \rightarrow$ give asymptotic but not exact value.

CONCLUSIONS AND FUTURE WORK

- We present the first distributional analysis of a linear probing hashing algorithm with buckets
- Key component of the analysis is the introduction of a new family of sequences $T_{k,d,b}$
- Better understanding of the properties of $T_{k,d,b}$, like combinatorial identities, asymptotic expansions. \Rightarrow complete understanding of occupancy distributions in buckets
- Combinatorial relation with other problems like Random Walks?
- Find other problems where $T_{k,d,b}$ may appear.