

# An optimal cardinality estimation algorithm based on order statistics and its full analysis

Jérémie Lumbroso

Projet ALGORITHMS, INRIA Rocquencourt &  
Laboratoire d'Informatique de Paris 6 (UPMC).

**AofA'10**

July 2, 2010—Vienna

# 1. PROBLEM STATEMENT

**Definition:** a stream  $\mathcal{S}$ ,

$$\mathcal{S} = s_1 s_2 \cdots s_N.$$

size of the stream:

$$|\mathcal{S}| = N$$

cardinality (= number of distinct elements):

$$\|\mathcal{S}\| = n$$

For instance,

$$\mathcal{S} = \text{run, sally, run, see, sally, run} \quad |\mathcal{S}| = 6 \quad \|\mathcal{S}\| = 3.$$

# 1. PROBLEM STATEMENT

**Definition:** a stream  $\mathcal{S}$ ,

$$\mathcal{S} = s_1 s_2 \cdots s_N.$$

size of the stream:

$$|\mathcal{S}| = N$$

cardinality (= **number of distinct elements**):

$$\|\mathcal{S}\| = n$$

For instance,

$$\mathcal{S} = \text{run, sally, run, see, sally, run} \quad |\mathcal{S}| = 6 \quad \|\mathcal{S}\| = 3.$$

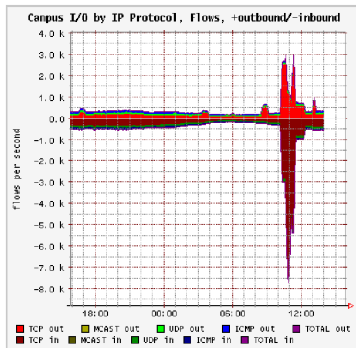
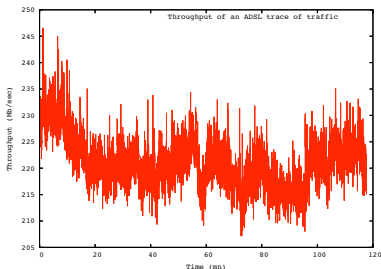
**Problem:** Estimate **cardinality** of  $\mathcal{S}$  so **large** that it **cannot be stored**.

## Constraints

- ▶ **very little** processing memory
- ▶ **on the fly** (single pass + simple main loop)
- ▶ **no** statistical hypothesis
- ▶ **accuracy** within a few percentiles

# MOTIVATION

- ▶ **Network security:**  
detect attacks (denial of service), or the spreading of worms/spam,...



- ▶ **Data mining:** document classification, ...
- ▶ **Databases:** query optimization
- ▶ **Distributed:** censor networks

## Bibliographic context

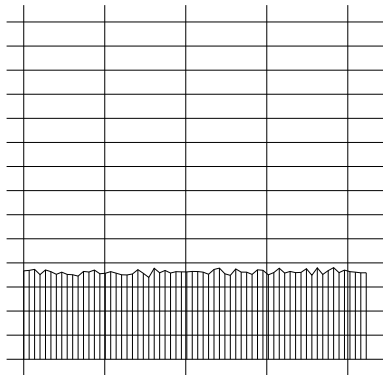
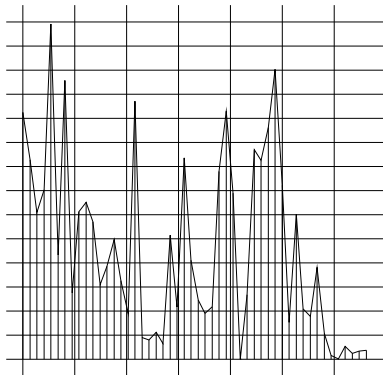
1. Algorithms based on pattern observation
  - ▶ Flajolet and Martin, 1985, *Probabilistic Counting*
  - ▶ Durand and Fl., 2003, *Loglog* ; Fl. and al., 2007, *Hyperloglog*
2. Algorithms based on order statistics
  - ▶ Giroire, 2003-2006, thèse P6
3. Complexity results
  - ▶ Alon, Matias, Szegedy, 1996, Frequency moments
  - ▶ Chassaing and Gerin, 2006, Theoretical optimality of using the minimum

## 2. THE MODEL

**Definition:** a hash function  $h$  is defined as

$$h : \mathcal{A}^* \rightarrow [0, 1].$$

**Main idea.** With “good enough” hash functions, our data is **uniformized**.



**Definition:** an **observable** = function of the underlying hash set  
(i.e.: a function **not sensitive to repetitions**)

Example: **minimum**

$$\blacktriangleright \min \{1, 2, 3\} = \min \{1, \dots, 1, 2, \dots, 2, 3, \dots, 3\} = 1$$

With:

- ▶ hash functions that **uniformize** the data
- ▶ **observables** : functions of underlying hash set

process data  $\rightarrow$  study  **$n$  i.i.d. random variables in  $[0, 1]$** .

### 3. ORDER STATISTIC of rank 1

$M$  := minimum of  $n$  random variables

$$\mathbb{P}_n[M \in [x, x + dx]] = n(1-x)^{n-1}dx \quad (1)$$

hence

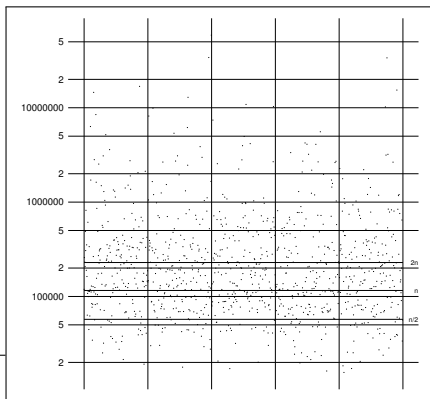
$$\mathbb{E}_n[M] = \int_0^1 x \cdot n(1-x)^{n-1}dx = \boxed{\frac{1}{n+1}}. \quad (2)$$

Advantages:

- ▶ computable in **one pass**
- ▶ computable with a **single register**
- ▶  $\mathbb{E}_n[M] = \frac{1}{n+1}$ .

Disadvantages:

- ▶ minimum **oscillates a lot**,  
indeed  $\sigma_n[M] = \frac{1}{n+1}$
- ▶ function  $x \mapsto \frac{1}{x}$  diverges at 0.



---

$$\mathbb{P}_n[M \leq \frac{t}{n}] \sim 1 - \exp(-t)$$



## STOCHASTIC AVERAGING

How to **cheaply** repeat the estimation (to average)?

**Idea.** Make  $m$  copies of the  $[0, 1]$  interval, and distribute the hashed values on these  $m$  intervals.

An extra condition: a **given element** must always be attributed to the **same interval**.

## Core Algorithm

Parameter:  $m$  control parameter

Input: a stream  $S = (s_1, \dots, s_N)$

---

**initialize**  $m$  registers  $M_1$  through  $M_m$  to 1

**forall**  $x \in S$  **do**

$A := h(x)$	{hash $x$ , with $h(x) \in (0, 1)$ }
$j := \lfloor mA \rfloor + 1$	{index of the substream assigned to $x$ }
$M_j := \min(M_j, mA - \lfloor mA \rfloor)$	{update minimum of $j$ -th substream}

**return**  $Z^* = m \cdot \frac{(m-1)}{M_1 + \dots + M_m}$

## 4. ANALYSIS of the CORE algorithm

$$\mathcal{Z}^* = m \cdot \frac{(m-1)}{M_1 + \dots + M_m}$$

Configuration  $\mathcal{C}$  defined by:

- ▶ allocation of  $n$  RVs in  $m$  bins (stochastic averaging)
- ▶ minimum of each bin

$$\mathbb{P}_n[\mathcal{C}] = \frac{1}{m^n} \binom{n}{n_1, \dots, n_m} \prod_{j=1}^m n_j (1 - x_j)^{n_j - 1} dx_j. \quad (3)$$

**[Interm.] Lemma.** The  $r$ -th moment of random variable  $\mathcal{Z}^*$  is given by

$$\mathbb{E}_n[(\mathcal{Z}^*)^r] = \bullet \int_{[0, \frac{n}{m}]^m} \left( 1 - \frac{1}{n} \sum_{j=1}^m t_j \right)^{n-m} \frac{dt_1 \cdots dt_m}{(t_1 + \dots + t_m)^r} \quad (4)$$

**Proof.**

1. sum (3) over all configurations
2. integrate over all possible minima:  $x_j \in [0, 1]$
3. rescaling ( $x_j = m/n \cdot t_j$ ) and algebraic manipulations.

## Calculating the multi-dimensional parametered integral

$$\mathbb{E}_n[(Z^*)^r] = \bullet \int_{[0, \frac{n}{m}]^m} \left(1 - \frac{1}{n} \sum_{j=1}^m t_j\right)^{n-m} \frac{dt_1 \cdots dt_m}{(t_1 + \dots + t_m)^r}$$

### Laplace method:

1. split integral into

$$I_C = \left[0, \frac{\delta(n)}{m}\right]^m \quad \text{and} \quad I_T = \left[0, \frac{n}{m}\right]^m \setminus I_C$$

2. on  $I_C$ , use  $(1 - \frac{1}{n} \sum) \sim \exp(-\sum)$
3. show  $I_T$  is negligible

+ use integral representation of Gamma function on integers

$$\int_0^\infty e^{-ay} a^{r-1} da = \frac{(r-1)!}{y^r}.$$

## A) UNBIASED and ACCURATE

**Theorem 1:**  $Z^*$  is asymptotically **unbiased**, in the sense that

$$\mathbb{E}_n[Z^*] = n(1 + o(1)). \quad (5)$$

**Theorem 2:** The precision of estimator  $Z^*$ , expressed in terms of standard error, satisfies

$$\frac{\sigma_n[Z^*]}{n} \sim \frac{1}{\sqrt{m-2}}. \quad (6)$$

## B) LIMIT DISTRIBUTION

Let  $S := M_1 + \dots + M_m$ , where the  $M_j$  are **interdependent**, the **Laplace transform**,

$$\mathbb{E}[e^{-wS}] \sim \left( \int_0^\infty e^{-t} e^{-w \frac{m}{n} t} dt \right)^m = (\mathbb{E}[e^{-wY \frac{m}{n}}])^m \quad (7)$$

and  $Y \in \text{Exp}(1)$ . So sum  $S$  behaves like **the sum of  $m$  indep.  $\text{Exp}(n/m)$** .

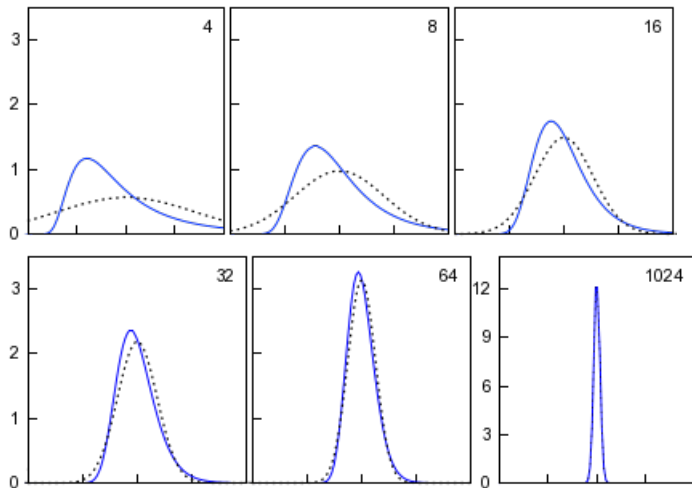
Thus, the rescaled/inverted  $\frac{1}{n/m \cdot S}$  has induced density:

$$\bar{w}_m(u) = e^{-1/u} \frac{u^{-m-1}}{(m-1)!}. \quad (8)$$

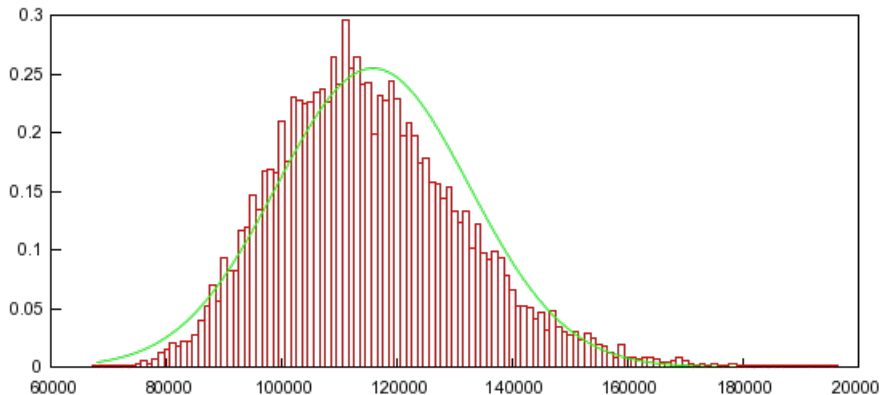
**Theorem 3:** For a fixed  $m > 1$ , as  $n$  tends to infinity, the estimator  $Z^*$  satisfies

$$\lim_{n \rightarrow \infty} \mathbb{P}_n \left[ \frac{Z^*}{n} \leq y \right] = \int_0^{y/(m-1)} e^{-1/u} \frac{u^{-m-1}}{(m-1)!} du.$$

# LIMIT DISTRIBUTION versus GAUSSIAN for $m = 4..1024$



## OBSERVED estimations for $m = 50$



Cardinality estimation for an English dictionary,  $n = 115\,794$  and  $m = 50$ .

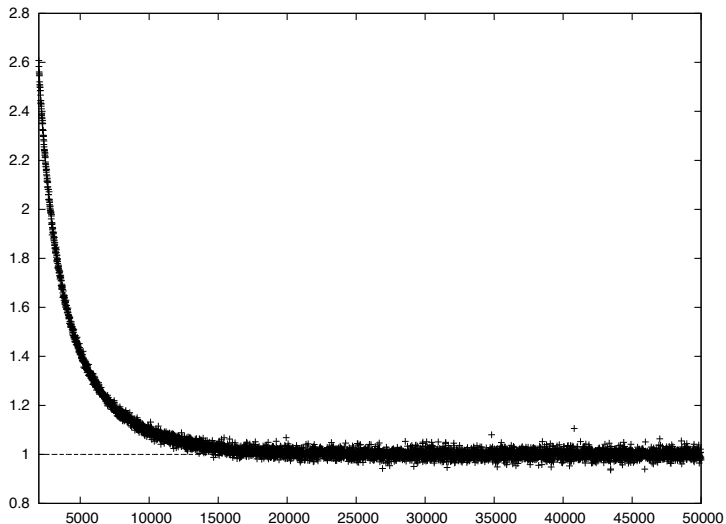
Gaussian with  $\mu = 115\,794$  and  $\sigma = 1/\sqrt{48}$ .



## 5. NON-ASYMPTOTICAL corrections

Plot of  $Z^*/n$  relative to  $n = 2000..50000$

(with  $m = 4096$ )



## A) POISSON model (SOME urns are empty)

**Pre-asymptotic calculations:**  $n$  relatively small compared to  $m$

Empty urns **bias** the average  $\rightarrow$  **keep track + ignore empty urns.**

Poisson approx. of urn allocation (i.e.:  $N_j \in \text{Poi}(\lambda)$ ).

$$\mathbb{P}[M_j \in [x, x + dx]] = \lambda e^{-\lambda x} dx + e^{-\lambda} \mathbf{1}_{\{x=1\}} \quad (9)$$

Let  $k = \#\{\text{non-empty urns}\}$ :

$$\mathbb{E} \left[ \frac{1}{M_1 + \dots + M_m} \right] = \sum_{k=0}^m \int_0^\infty \binom{n}{k} \left( \lambda \cdot \frac{1 - e^{-a-\lambda}}{a + \lambda} \right)^k (e^{-a-\lambda})^{m-k} da$$

after calculations + Laplace, yields

**Theorem 4:** let  $n/m = \lambda$  be such that  $0 < \lambda < C < \infty$ , then

$$\mathbb{E}_n[\mathcal{Z}^*] \sim \frac{n}{1 - e^{-\lambda}}. \quad (10)$$

## B) “Linear counting”<sup>1</sup> (TOO MANY urns are empty)

**Non-asymptotic:** shift in point-of-view

$n$  balls are thrown into  $m$  urns.

**Classical result:** Let  $W_k := \#\{\text{urns containing } k \text{ balls}\}$

$$\mathbb{E}[W_k] \sim m \cdot \left[ \frac{\lambda^k}{k!} \exp(-\lambda) \right]$$

where  $\lambda := n/m$  is constant, with  $n \rightarrow \infty$  and  $m \rightarrow \infty$ .

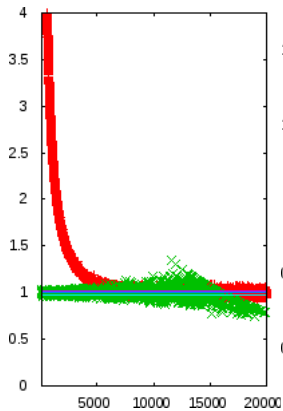
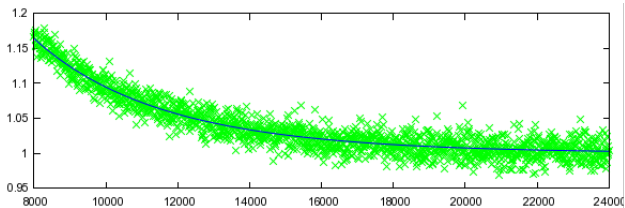
Since  $\mathbb{E}[W_0] \sim m \cdot \exp\left(-\frac{n}{m}\right)$  then, with  $\hat{w}_0$  observed empty urns,

$$n \approx -m \log \left( \frac{\hat{w}_0}{m} \right) \quad (11)$$

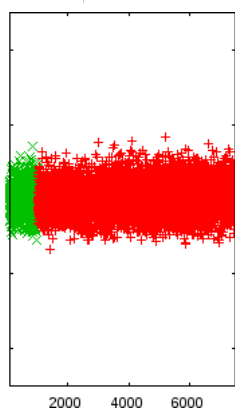
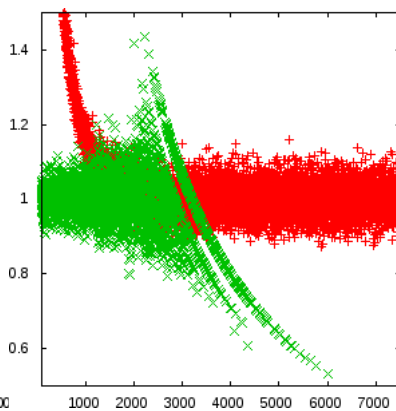
---

<sup>1</sup>after Whang et al, “A Linear-Time Probabilistic Counting Algorithm for Database Applications,” ACM Trans. on Database Systems, Vol. 15, No. 2, pp. 208-229, June 1990.

## C) JOINING all regimes



linear counting



core algorithm 19/20

## 5. CONCLUSION

- ▶ a **complete** algorithm (large range of cardinalities)
- ▶ **optimal** within its class
- ▶ the **full** analysis with precise knowledge of the **output**

## 5. CONCLUSION

- ▶ a **complete** algorithm (large range of cardinalities)
- ▶ **optimal** within its class
- ▶ the **full** analysis with precise knowledge of the **output**

Future (immediate):

- ▶ **rate of convergence**
- ▶ attempt transposing the analysis to Hyperloglog
- ▶ plug-in the limit distribution analysis in simple algorithms which use cardinality estimation as a black box.