# The Analysis of Data Stream Algorithms

Conrado Martínez
Universitat Politècnica de Catalunya

AofA2023@Taipei, June 26–30, 2023

# Outline of the Talk

# Part I

## Random Musings

# The fate of AofA (in CS)?

# The fate of AofA (in CS)?

While many of the techniques and results of our area are received with interest and recognition outside Computer Science, is it the case for CS anymore?

Do our papers often succeed in major CS journals and conferences?

Do our results have a noticeable impact in other CS research communities?

# The fate of AofA (in CS)?

Troubles in AofA land:

- Are we analyzing algorithms and data structures often enough (aren't we AofA?)
- Are we trying hard enough to provide new insights and useful answers in many striving areas of CS?
- Drawback: scientific mindset (why does it work? why doesn't it work!?) vs engineering mindset (does it work? how does it work?)

# The fate of AofA (in CS)?

A roadmap to find the way in AofA land:

1. Consider submitting your work to major (T)CS conferences and journals

2. Advocate for AofA methods and results whenever the opportunity arises

3. Promote the scientific mindset in CS (not just in the theoretical areas!) $\rightarrow$ Algorithm Science

4. Look for problems where a precise probabilistic analysis is crucial or the only reasonable option

5. "Package" your results in general form (maybe as software tools?) and try to make them easy to use and to benefit from

# The fate of AofA (in CS)?

Promising lands for AofA, some new, some old:

1. Data streaming algorithms
2. Similarity & proximity search
3. Randomized metaheuristics: EA, GAs, Simulated annealing, ACO, PSO, ... ($\rightarrow$ Benjamin's talk)
4. Deep learning & stochastic gradient descent ($\rightarrow$ Chih-Jen's talk)
5. Data & process mining
6. ...

# A personal account



Ph. Flajolet

My first incursion into the very rich area of data stream algorithms dates back to 2011, and I am still interested and in love with it. Isn't it ironic that 2011 was the year that Flajolet passed away 😢😢?

Philippe, the person who, besides many other fundamental achievements in AofA & Analytic Combinatorics, had developed some of the most elegant and practical algorithms in the area, beginning with his celebrated Probabilistic Count together with G.N. Martin in the mid eighties.

So let's move on ... and talk a bit about Data Stream Algorithms and the fundamental contributions of AofA to the area!



**Philippe Flajolet**
Father of Approximate Counting Algorithms

I can count in less memory if you do not want exact answer.

# Part II

## Introduction

# Introduction

- A data stream is a (very long) sequence

$$\mathcal{Z} = z_1, z_2, z_3, \ldots, z_N$$

of elements drawn from a (very large) domain $\mathcal{U}$ ($z_i \in \mathcal{U}$)

- The goal: to compute $f(\mathcal{Z})$, but ...

# Introduction

- A data stream is a (very long) sequence

$$\mathcal{Z} = z_1, z_2, z_3, \ldots, z_N$$

of elements drawn from a (very large) domain $\mathcal{U}$ ($z_i \in \mathcal{U}$)
- The goal: to compute $f(\mathcal{Z})$, but ...

# Introduction

... under rather stringent constraints (data stream model)

- a single pass over the data stream
- extremely short time spent on each single data item
- a limited amount $M$ of auxiliary memory, $M \ll N$; ideally $M = \Theta(1)$ or $M = \Theta(\log N)$
- no statistical hypothesis about the data

# Introduction

. . . under rather stringent constraints (data stream model)

- a single pass over the data stream
- extremely short time spent on each single data item
- a limited amount $M$ of auxiliary memory, $M \ll N$; ideally $M = \Theta(1)$ or $M = \Theta(\log N)$
- no statistical hypothesis about the data

# Introduction

. . . under rather stringent constraints (data stream model)

- a single pass over the data stream
- extremely short time spent on each single data item
- a limited amount $M$ of auxiliary memory, $M \ll N$; ideally $M = \Theta(1)$ or $M = \Theta(\log N)$
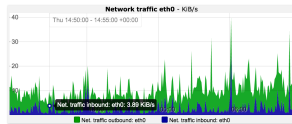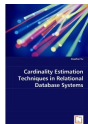- no statistical hypothesis about the data

# Introduction

... under rather stringent constraints (data stream model)

- a single pass over the data stream
- extremely short time spent on each single data item
- a limited amount $M$ of auxiliary memory, $M \ll N$; ideally $M = \Theta(1)$ or $M = \Theta(\log N)$
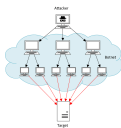- no statistical hypothesis about the data

There is a wide range of applications for the data stream model

- Network traffic analysis $\Rightarrow$ DoS/DDoS attacks, *worms*, . . .
- Database query optimization
- Information retrieval $\Rightarrow$ similarity index
- Data mining
- Recommedation systems
- and many more . . .

# Introduction



There is a wide range of applications for the data stream model

- Network traffic analysis $\Rightarrow$ DoS/DDoS attacks, *worms*, ...
- Database query optimization
- Information retrieval $\Rightarrow$ similarity index
- Data mining
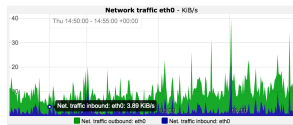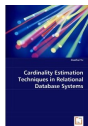- Recommedation systems
- and many more ...

# Introduction



There is a wide range of applications for the data stream model

- Network traffic analysis $\Rightarrow$ DoS/DDoS attacks, *worms*, . . .
- Database query optimization
- Information retrieval $\Rightarrow$ similarity index
- Data mining
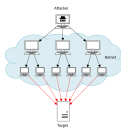- Recommedation systems
- and many more . . .

# Introduction



There is a wide range of applications for the data stream model

- Network traffic analysis $\Rightarrow$ DoS/DDoS attacks, *worms*, ...
- Database query optimization
- Information retrieval $\Rightarrow$ similarity index
- Data mining
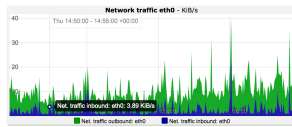- Recommedation systems
- and many more ...

# Introduction



There is a wide range of applications for the data stream model

- Network traffic analysis $\Rightarrow$ DoS/DDoS attacks, *worms*, . . .
- Database query optimization
- Information retrieval $\Rightarrow$ similarity index
- Data mining
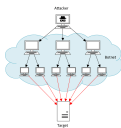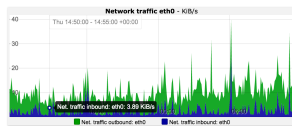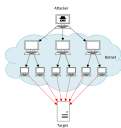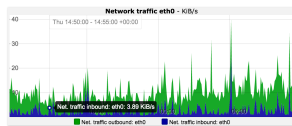- Recommedation systems
- and many more . . .

# Introduction



There is a wide range of applications for the data stream model

- Network traffic analysis $\Rightarrow$ DoS/DDoS attacks, *worms*, ...
- Database query optimization
- Information retrieval $\Rightarrow$ similarity index
- Data mining
- Recommedation systems
- and many more ...

# Introduction

We'll often look at $\mathcal{Z}$ as a multiset $\{x_1 \circ f_1, \ldots, x_n \circ f_n\}$, where

$$f_i = \text{frequency of the } i\text{-th distinct element } x_i$$

Some fundamental problems in data stream analysis:

- <u>Number of distinct elements</u>: $\text{card}(\mathcal{Z}) = n \leqslant N$
- Random samples of distinct elements
- Frequency moments $F_p = \sum_{1 \leqslant i \leqslant n} f_i^p$
  (N.B. $n = F_0$, $N = F_1$)
- (Number of) Elements $x_i$ such that $f_i \geqslant k$ (k-elephants) or $f_i < k$ (k-mice)
- (Number of) Elements $x_i$ such that $f_i/N \geqslant c$, $0 < c < 1$ (c-icebergs, a.k.a. heavy hitters)
- The k most frequent elements (top-k elements)

# Introduction

We'll often look at $\mathcal{Z}$ as a multiset $\{x_1 \circ f_1, \dots, x_n \circ f_n\}$, where

$$f_i = \text{frequency of the } i\text{-th distinct element } x_i$$

Some fundamental problems in data stream analysis:

- <u>Number of distinct elements</u>: $\text{card}(\mathcal{Z}) = n \leqslant N$
- <u>Random samples of distinct elements</u>
- Frequency moments $F_p = \sum_{1 \leqslant i \leqslant n} f_i^p$
  (N.B. $n = F_0, N = F_1$)
- (Number of) Elements $x_i$ such that $f_i \geqslant k$ (k-elephants) or $f_i < k$ (k-mice)
- (Number of) Elements $x_i$ such that $f_i/N \geqslant c$, $0 < c < 1$ (c-icebergs, a.k.a. heavy hitters)
- The k most frequent elements (top-k elements)

# Introduction

We'll often look at $\mathcal{Z}$ as a multiset $\{x_1 \circ f_1, \ldots, x_n \circ f_n\}$, where

$$f_i = \text{frequency of the } i\text{-th distinct element } x_i$$

Some fundamental problems in data stream analysis:

- <u>Number of distinct elements</u>: $\text{card}(\mathcal{Z}) = n \leqslant N$
- <u>Random samples of distinct elements</u>
- Frequency moments $F_p = \sum_{1 \leqslant i \leqslant n} f_i^p$
  (N.B. $n = F_0, N = F_1$)
- (Number of) Elements $x_i$ such that $f_i \geqslant k$ (k-elephants) or $f_i < k$ (k-mice)
- (Number of) Elements $x_i$ such that $f_i/N \geqslant c$, $0 < c < 1$ (c-icebergs, a.k.a. heavy hitters)
- The $k$ most frequent elements (top-k elements)

# Introduction

We'll often look at $\mathcal{Z}$ as a multiset $\{x_1 \circ f_1, \ldots, x_n \circ f_n\}$, where

$$f_i = \text{frequency of the } i\text{-th distinct element } x_i$$

Some fundamental problems in data stream analysis:

- <u>Number of distinct elements</u>: $\text{card}(\mathcal{Z}) = n \leqslant N$
- <u>Random samples of distinct elements</u>
- <u>Frequency moments</u> $F_p = \sum_{1 \leqslant i \leqslant n} f_i^p$
  (N.B. $n = F_0$, $N = F_1$)
- (Number of) Elements $x_i$ such that $f_i \geqslant k$ (k-elephants) or $f_i < k$ (k-mice)
- (Number of) Elements $x_i$ such that $f_i/N \geqslant c$, $0 < c < 1$ (c-icebergs, a.k.a. heavy hitters)
- The $k$ most frequent elements (top-k elements)

# Introduction

We'll often look at $\mathcal{Z}$ as a multiset $\{x_1 \circ f_1, \ldots, x_n \circ f_n\}$, where

$$f_i = \text{frequency of the } i\text{-th distinct element } x_i$$

Some fundamental problems in data stream analysis:

- <u>Number of distinct elements</u>: $\mathrm{card}(\mathcal{Z}) = n \leqslant N$
- <u>Random samples of distinct elements</u>
- Frequency moments $F_p = \sum_{1 \leqslant i \leqslant n} f_i^p$
  (N.B. $n = F_0, N = F_1$)
- (Number of) Elements $x_i$ such that $f_i \geqslant k$ (k-elephants) or $f_i < k$ (k-mice)
- (Number of) Elements $x_i$ such that $f_i/N \geqslant c$, $0 < c < 1$
  (c-icebergs, a.k.a. heavy hitters)
- The $k$ most frequent elements (top-k elements)

# Introduction

We'll often look at $\mathcal{Z}$ as a multiset $\{x_1 \circ f_1, \ldots, x_n \circ f_n\}$, where

$$f_i = \text{frequency of the } i\text{-th distinct element } x_i$$

Some fundamental problems in data stream analysis:

- <u>Number of distinct elements</u>: $\text{card}(\mathcal{Z}) = n \leqslant N$
- <u>Random samples of distinct elements</u>
- Frequency moments $F_p = \sum_{1 \leqslant i \leqslant n} f_i^p$
  (N.B. $n = F_0, N = F_1$)
- (Number of) Elements $x_i$ such that $f_i \geqslant k$ (k-elephants) or $f_i < k$ (k-mice)
- (Number of) Elements $x_i$ such that $f_i/N \geqslant c$, $0 < c < 1$ (c-icebergs, a.k.a. heavy hitters)
- The $k$ most frequent elements (top-k elements)

# Introduction

Very limited available memory $\Rightarrow$ exact solution too costly or unfeasible

$\Rightarrow$ Randomized algorithms $\Rightarrow$ estimation $\hat{q}$ of the quantity of interest $q = f(\mathcal{Z})$

- $\hat{q}$ must be an unbiased estimator

$$E[\hat{q}] = q$$

- The estimator must accurate, for example, it must have a small standard error

$$SE[\hat{q}] := \frac{\sqrt{Var[\hat{q}]}}{E[\hat{q}]} < \epsilon,$$

e.g., $\epsilon = 0.01$ (1%)

# Introduction

Very limited available memory $\Rightarrow$ exact solution too costly or unfeasible
$\Rightarrow$ Randomized algorithms $\Rightarrow$ estimation $\hat{q}$ of the quantity of interest $q = f(\mathcal{Z})$

- $\hat{q}$ must be an unbiased estimator

$$E[\hat{q}] = q$$

- The estimator must accurate, for example, it must have a small standard error

$$SE[\hat{q}] := \frac{\sqrt{Var[\hat{q}]}}{E[\hat{q}]} < \epsilon,$$

e.g., $\epsilon = 0.01$ (1%)

# Part III

## Cardinality Estimation

**1** Probabilistic Counting

**2** LogLog & HyperLogLog

**3** Order Statistics

**4** Recordinality

# Probabilistic Counting



G.N. Martin

In late 70s G. Nigel Martin invented probabilistic counting to optimize database query performance

To correct the bias that he systematically found in his experiments, he introduced a "fudge" factor in the estimator

# Probabilistic Counting

When Philippe Flajolet learnt about the algorithm, he put it on a solid scientific ground, with a detailed mathematical analysis which delivered the exact value of the correction factor and a tight upper bound on the standard error

As I said over the phone, I started working on your algorithm when Kyu-Young Whang considered implementing it and wanted explanations/estimations. I find it simple, eleg and amazingly powerful.

# Probabilistic Counting

- **Key idea**: every element is hashed to a real value in $(0, 1)$ $\Rightarrow$ reproducible randomness
- The "multiset" $\mathcal{Z}$ is mapped by the hash function $h : \mathcal{U} \to (0, 1)$ to a multiset

$$\mathcal{Z}' = h(\mathcal{Z}) = \{y_1 \circ f_1, \ldots, y_n \circ f_n\},$$

with $y_i = \text{hash}(x_i)$, $f_i = $ frequency of $x_i$ in $\mathcal{Z}$
- The set of distinct* elements $Y = \{y_1, \ldots, y_n\}$ is a set of $n$ random numbers, independent and uniformly drawn from $(0, 1)$

# Probabilistic Counting

- **Key idea**: every element is hashed to a real value in $(0, 1) \Rightarrow$ reproducible randomness
- The "multiset" $\mathcal{Z}$ is mapped by the hash function $h : \mathcal{U} \to (0, 1)$ to a multiset

$$\mathcal{Z}' = h(\mathcal{Z}) = \{y_1 \circ f_1, \ldots, y_n \circ f_n\},$$

with $y_i = \mathsf{hash}(x_i)$, $f_i = $ frequency of $x_i$ in $\mathcal{Z}$

- The set of distinct* elements $Y = \{y_1, \ldots, y_n\}$ is a set of $n$ random numbers, independent and uniformly drawn from $(0, 1)$

# Probabilistic Counting

- **Key idea**: every element is hashed to a real value in $(0, 1) \Rightarrow$ reproductible randomness
- The "multiset" $\mathcal{Z}$ is mapped by the hash function $h : \mathcal{U} \to (0, 1)$ to a multiset

$$\mathcal{Z}' = h(\mathcal{Z}) = \{y_1 \circ f_1, \ldots, y_n \circ f_n\},$$

  with $y_i = \text{hash}(x_i)$, $f_i = $ frequency of $x_i$ in $\mathcal{Z}$
- The set of distinct* elements $Y = \{y_1, \ldots, y_n\}$ is a set of $n$ random numbers, independent and uniformly drawn from $(0, 1)$

*We'll neglect the probability of collisions, i.e., $h(x_i) = h(x_j)$ for some $x_i \neq x_j$; this is reasonable if $h(x)$ has enough bits

# Probabilistic Counting

Flajolet & Martin (JCSS, 1985) proposed to find, among the set of hash values, the length of the largest R such that hash values with the prefix $0.0^{p-1}1\ldots$, have appeared in the stream, for all $p$, $1 \leqslant p \leqslant R$

The value R is an observable which can be easily be computed using a small auxiliary memory and it is insensitive to repetitions ← the observable is a function of Y, not of the $f_i$'s

# Probabilistic Counting

- For a set of $n$ random numbers in $(0, 1) \rightarrow$

$$E\left[R\right] \approx \log_2 n$$

- However $E\left[2^R\right] \not\sim n$, there is a significant bias and we need $\phi$ such that

$$E\left[\phi \cdot 2^R\right] \sim n$$

# Probabilistic Counting

- For a set of $n$ random numbers in $(0, 1) \rightarrow$

$$\mathsf{E}\left[R\right] \approx \log_2 n$$

- However $\mathsf{E}\left[2^R\right] \not\sim n$, there is a significant bias and we need $\phi$ such that

$$\mathsf{E}\left[\phi \cdot 2^R\right] \sim n$$

# Probabilistic Counting

```
procedure ProbabilisticCounting(Z)
    bmap ← ⟨0, 0, . . . , 0⟩
    for z ∈ Z do
        y ← hash(z)
        p ← lenght of the largest prefix 0.0^{p-1}1... in y
        bmap[p] ← 1
    end for
    R ← largest p such that bmap[i] = 1 for all 1 ⩽ i ⩽ p
    // φ is the correction factor: E[φ · 2^R] = n
    return  Z := φ · 2^R
end procedure
```

A very precise mathematical analysis gives
($\gamma$ = Euler's gamma constant,
$\nu(k)$ = # of 1's in binary repr. of k):

$$\phi^{-1} = \frac{e^{\gamma}\sqrt{2}}{3} \prod_{k \geqslant 1} \left( \frac{(4k+1)(2k+1)}{2k(4k+3)} \right)^{(-1)^{\nu(k)}} \approx 0.77351\ldots$$

# Stochastic averaging

- The standard error of $Z := \phi \cdot 2^R$, despite constant, is too large: $SE[Z] > 1$
- Second key idea: "repeat" several times to reduce variance and improve precision
- Problem: using $m$ hash functions to generate $m$ streams is too costly and it's very difficult to guarantee independence between the hash values

# Stochastic averaging

- The standard error of $Z := \phi \cdot 2^R$, despite constant, is too large: $SE[Z] > 1$
- Second key idea: "repeat" several times to reduce variance and improve precision
- Problem: using $m$ hash functions to generate $m$ streams is too costly and it's very difficult to guarantee independence between the hash values

# Stochastic averaging

- The standard error of $Z := \phi \cdot 2^R$, despite constant, is too large: $SE[Z] > 1$
- Second key idea: "repeat" several times to reduce variance and improve precision
- Problem: using $\mathfrak{m}$ hash functions to generate $\mathfrak{m}$ streams is too costly and it's very difficult to guarantee independence between the hash values

# Stochastic averaging



- Use the first $\log_2 m$ bits of each hash value to "redirect" it (the remaining bits) to one of the $m$ substreams $\rightarrow$ stochastic averaging

- Obtain $m$ observables $R_1, R_2, \ldots, R_m$, one from each substream

- Each $R_i$ can give us an estimation for the cardinality of the $i$-th substream, namely, $R_i$ can be used to estimate $n/m$; the mean value $\bar{R} = 1/m \sum R_i$ can also be used to estimate $n/m$

# Stochastic averaging



- Use the first $\log_2 m$ bits of each hash value to "redirect" it (the remaining bits) to one of the $m$ substreams $\rightarrow$ stochastic averaging
- Obtain $m$ observables $R_1, R_2, \ldots, R_m$, one from each substream
- Each $R_i$ can give us an estimation for the cardinality of the $i$-th substream, namely, $R_i$ can be used to estimate $n/m$; the mean value $\bar{R} = 1/m \sum R_i$ can also be used to estimate $n/m$

# Stochastic averaging



- Use the first $\log_2 m$ bits of each hash value to "redirect" it (the remaining bits) to one of the $m$ substreams $\rightarrow$ stochastic averaging
- Obtain $m$ observables $R_1, R_2, \ldots, R_m$, one from each substream
- Each $R_i$ can give us an estimation for the cardinality of the $i$-th substream, namely, $R_i$ can be used to estimate $n/m$; the mean value $\overline{R} = 1/m \sum R_i$ can also be used to estimate $n/m$

# Stochastic averaging

There are many different options to compute an estimator from the $m$ observables

- Sum of estimators:

$$Z_1 := \phi_1(2^{R_1} + \ldots + 2^{R_m})$$

- Arithmetic mean of observables (as proposed by Flajolet & Martin):

$$Z_2 := m \cdot \phi_2 \cdot 2^{\frac{1}{m} \sum_{1 \leqslant i \leqslant m} R_i}$$

# Stochastic averaging

- **Harmonic mean** (keep tuned):

$$Z_3 := \phi_3 \cdot \frac{m^2}{2^{-R_1} + 2^{-R_2} + \ldots + 2^{-R_m}}$$

Since $2^{-R_i} \approx m/n$, the second factor gives $\approx m^2/(m^2/n) = n$

# Stochastic averaging

- All the strategies above yield a standard error of the form

$$\frac{c}{\sqrt{m}} + \text{l.o.t.}$$

Larger memory $\Rightarrow$ improved precision!

- In *probabilistic counting* the authors used the arithmetic mean of observables

$$SE\left[Z_{\text{ProbCount}}\right] \approx \frac{0.78}{\sqrt{m}}$$

# Stochastic averaging

- All the strategies above yield a standard error of the form

$$\frac{c}{\sqrt{m}} + \text{l.o.t.}$$

  Larger memory $\Rightarrow$ improved precision!

- In *probabilistic counting* the authors used the arithmetic mean of observables

$$SE\left[Z_{\mathsf{ProbCount}}\right] \approx \frac{0.78}{\sqrt{m}}$$

# LogLog & HyperLogLog



M. Durand

- Durand & Flajolet (2003) realized that the bitmaps ($\Theta(\log n)$ bits) used by *Probabilistic Counting* can be avoided and propose as observable the largest R such that the pattern $0.0^{R-1}1$ appears

- The new observable is similar to that of *Probabilistic Counting* but not equal: $R(LogLog) \geqslant R(ProbCount)$

  Example

  Observed patterns: $0.1101\ldots$, $0.010\ldots$, $0.0011\ldots$, $0.00001\ldots$

  $R(LogLog) = 5$, $R(ProbCount) = 3$

# LogLog & HyperLogLog



M. Durand

- Durand & Flajolet (2003) realized that the bitmaps ($\Theta(\log n)$ bits) used by *Probabilistic Counting* can be avoided and propose as observable the largest R such that the pattern $0.0^{R-1}1$ appears
- The new observable is similar to that of *Probabilistic Counting* but not equal: $R(\text{LogLog}) \geqslant R(\text{ProbCount})$

> **Example**
>
> Observed patterns: $0.1101\ldots$, $0.010\ldots$, $0.0011\ldots$, $0.00001\ldots$
>
> $R(\text{LogLog}) = 5,$ $\qquad R(\text{ProbCount}) = 3$

# LogLog & HyperLogLog



M. Durand

- Durand & Flajolet (2003) realized that the bitmaps ($\Theta(\log n)$ bits) used by *Probabilistic Counting* can be avoided and propose as observable the largest R such that the pattern $0.0^{R-1}1$ appears
- The new observable is similar to that of *Probabilistic Counting* but not equal: $R(\text{LogLog}) \geqslant R(\text{ProbCount})$

> **Example**
>
> Observed patterns: $0.1101\ldots$, $0.010\ldots$, $0.0011\ldots$, $0.00001\ldots$
>
> $R(\text{LogLog}) = 5$, $\qquad R(\text{ProbCount}) = 3$

# LogLog & HyperLogLog

- The new observable is simpler to obtain: keep updated the largest R seen so far: $R := \max\{R, p\} \Rightarrow$ only $\Theta(\log \log n)$ bits needed, since $E[R] = \Theta(\log n)$!

- We have $E[R] \sim \log_2 n$, but $E\left[2^R\right] = +\infty$, *stochastic averaging* comes to rescue!

- For LogLog, Durand & Flajolet propose

$$Z_{\text{LogLog}} := \alpha_m \cdot m \cdot 2^{\frac{1}{m} \sum_{1 \leqslant i \leqslant m} R_i}$$

# LogLog & HyperLogLog

- The new observable is simpler to obtain: keep updated the largest R seen so far: $R := \max\{R, p\} \Rightarrow$ only $\Theta(\log \log n)$ bits needed, since $E[R] = \Theta(\log n)$!

- We have $E[R] \sim \log_2 n$, but $E[2^R] = +\infty$, *stochastic averaging* comes to rescue!

- For LogLog, Durand & Flajolet propose

$$Z_{\mathsf{LogLog}} := \alpha_m \cdot m \cdot 2^{\frac{1}{m} \sum_{1 \leqslant i \leqslant m} R_i}$$

# LogLog & HyperLogLog

- The new observable is simpler to obtain: keep updated the largest R seen so far: $R := \max\{R, p\} \Rightarrow$ only $\Theta(\log \log n)$ bits needed, since $E[R] = \Theta(\log n)$!
- We have $E[R] \sim \log_2 n$, but $E[2^R] = +\infty$, *stochastic averaging* comes to rescue!
- For LogLog, Durand & Flajolet propose

$$Z_{\mathsf{LogLog}} := \alpha_m \cdot m \cdot 2^{\frac{1}{m} \sum_{1 \leqslant i \leqslant m} R_i}$$

# LogLog & HyperLogLog

- The mathematical analysis gives for the correcting factor

$$\alpha_m = \left( \Gamma(-1/m) \frac{1 - 2^{1/m}}{\ln 2} \right)^{-m}$$

that guarantees that $E[Z] = n + \text{l.o.t.}$ (asymptotically unbiased) and the standard error is

$$SE\left[ Z_{\text{LogLog}} \right] \approx \frac{1.30}{\sqrt{m}}$$

- Only $m$ counters of size $\log_2 \log_2(n/m)$ bits needed: Ex.: $m = 2048 = 2^{11}$ counters, 5 bits each (1.25 Kbyte in total), are enough to give precise cardinality estimations for $n$ up to $2^{27} \approx 10^8$, with an standard error less than 4%

# LogLog & HyperLogLog

- The mathematical analysis gives for the correcting factor

$$\alpha_m = \left( \Gamma(-1/m) \frac{1 - 2^{1/m}}{\ln 2} \right)^{-m}$$

  that guarantees that $E[Z] = n + \text{l.o.t.}$ (asymptotically unbiased) and the standard error is
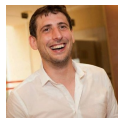
$$SE\left[ Z_{\text{LogLog}} \right] \approx \frac{1.30}{\sqrt{m}}$$

- Only $m$ counters of size $\log_2 \log_2(n/m)$ bits needed: Ex.: $m = 2048 = 2^{11}$ counters, 5 bits each (1.25 Kbyte in total), are enough to give precise cardinality estimations for $n$ up to $2^{27} \approx 10^8$, with an standard error less than 4%

# LogLog & HyperLogLog



É. Fusy    O. Gandouet    F. Meunier

- Flajolet, Fusy, Gandouet & Meunier conceived in 2007 the best algorithm known (cif. Flajolet's *keynote speech* in ITC Paris 2009)

- Briefly: HyperLogLog combines the LogLog observables $R_i$ using the harmonic mean instead of the arithmetic mean

$$SE\left[Z_{\text{HyperLogLog}}\right] \approx \frac{1.03}{\sqrt{m}}$$

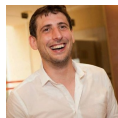# LogLog & HyperLogLog



É. Fusy     O. Gandouet     F. Meunier

- Flajolet, Fusy, Gandouet & Meunier conceived in 2007 the best algorithm known (cif. Flajolet's *keynote speech* in ITC Paris 2009)
- Briefly: HyperLogLog combines the LogLog observables $R_i$ using the harmonic mean instead of the arithmetic mean

$$SE\left[Z_{\text{HyperLogLog}}\right] \approx \frac{1.03}{\sqrt{m}}$$

# LogLog & HyperLogLog



P. Chassaing    L. Gerin

- The idea of HyperLogLog stems from the analytical study of Chassaing & Gerin (2006) to show the optimal way to combine observables, but in their study the observables were the $k$-th order statistics of each substream (next!)
- They proved that the optimal way to combine them is to use the harmonic mean

# LogLog & HyperLogLog



P. Chassaing    L. Gerin

- The idea of HyperLogLog stems from the analytical study of Chassaing & Gerin (2006) to show the optimal way to combine observables, but in their study the observables were the $k$-th order statistics of each substream (next!)
- They proved that the optimal way to combine them is to use the harmonic mean

# Order Statistics

- Bar-Yossef, Kumar & Sivakumar (2002); Bar-Yossef, Jayram, Kumar, Sivakumar & Trevisan (2002) have proposed to use the $k$-th order statistic $Y_{(k)}$ to estimate cardinality (KMV algorithm); for a set of $n$ random numbers, independent and uniformly distributed in $(0, 1)$

$$E\left[Y_{(k)}\right] = \frac{k}{n+1} \Rightarrow E\left[\frac{k-1}{Y_{(k)}}\right] = n$$

- Giroire (2005, 2009) also proposes several estimators combining order statistics via *stochastic averaging*

# Order Statistics

- Bar-Yossef, Kumar & Sivakumar (2002); Bar-Yossef, Jayram, Kumar, Sivakumar & Trevisan (2002) have proposed to use the $k$-th order statistic $Y_{(k)}$ to estimate cardinality (KMV algorithm); for a set of $n$ random numbers, independent and uniformly distributed in $(0, 1)$

$$E\left[Y_{(k)}\right] = \frac{k}{n+1} \Rightarrow E\left[\frac{k-1}{Y_{(k)}}\right] = n$$

- Giroire (2005, 2009) also proposes several estimators combining order statistics via *stochastic averaging*
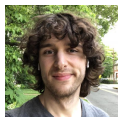
# Order Statistics



J. Lumbroso

- The minimum of the set ($k = 1$) does not allow a feasible estimator, but again *stochastic averaging* comes to rescue

- Lumbroso uses the mean of $m$ minima, one for each substream

$$Z_{\text{MinCount}} := \frac{m(m-1)}{M_1 + \ldots + M_m},$$

where $M_i$ is the minimum hash value of the $i$-th substream
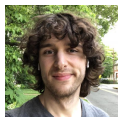
# Order Statistics



J. Lumbroso

- The minimum of the set ($k = 1$) does not allow a feasible estimator, but again *stochastic averaging* comes to rescue
- Lumbroso uses the mean of $m$ minima, one for each substream

$$Z_{\text{MinCount}} := \frac{m(m-1)}{M_1 + \ldots + M_m},$$

where $M_i$ is the minimum hash value of the $i$-th substream

# Order Statistics



- MinCount is an unbiased estimator with standard error $1/\sqrt{m-2}$
- Lumbroso also succeeds to compute the probability distribution of $Z_{MinCount}$ and the small corrections needed to estimate small cardinalities (too few elements hashing to one particular substream)

# Order Statistics



- MinCount is an unbiased estimator with standard error $1/\sqrt{m-2}$
- Lumbroso also succeeds to compute the probability distribution of $Z_{\mathrm{MinCount}}$ and the small corrections needed to estimate small cardinalities (too few elements hashing to one particular substream)

# Recordinality



A. Helmi    A. Viola

- Recordinality (Helmi, Lumbroso, M., Viola, 2012) is loosely related to order statistics, but based in completely different principles and it exhibits several unique features
- Some of the ideas where very useful to develop Affirmative Sampling, stay tuned!

# Recordinality



A. Helmi     A. Viola

- Recordinality (Helmi, Lumbroso, M., Viola, 2012) is loosely related to order statistics, but based in completely different principles and it exhibits several unique features
- Some of the ideas where very useful to develop Affirmative Sampling, stay tuned!

# Recordinality

- Recordinality counts the number of records (more generally, k-records) in the sequence of hash values
- It depends in the underlying permutation of the first occurrences of distinct values, very different from the other estimators
- If we assume that the first occurrences of distinct values form a random permutation then there's no need for hash values!

# Recordinality

- Recordinality counts the number of records (more generally, $k$-records) in the sequence of hash values
- It depends in the underlying permutation of the first occurrences of distinct values, very different from the other estimators
- If we assume that the first occurrences of distinct values form a random permutation then there's no need for hash values!

# Recordinality

- Recordinality counts the number of records (more generally, $k$-records) in the sequence of hash values
- It depends in the underlying permutation of the first occurrences of distinct values, very different from the other estimators
- If we assume that the first occurrences of distinct values form a random permutation then there's no need for hash values!

# Recordinality

- $\sigma(i)$ is a record of the permutation $\sigma$ if $\sigma(i) > \sigma(j)$ for all $j < i$
- This notion is generalized to k-records: $\sigma(i)$ is a k-record if there are at most $k-1$ elements $\sigma(j)$ larger than $\sigma(i)$ for $j < i$; in other words, $\sigma(i)$ is among the $k$ largest elements in $\sigma(1), \ldots, \sigma(i)$

  Example
  This example permutation contains six 2-records

  $$\sigma = 3, 6, 1, 12, 8, 10, 4, 13, 7, 5, 9, 11, 2$$

# Recordinality

- $\sigma(i)$ is a record of the permutation $\sigma$ if $\sigma(i) > \sigma(j)$ for all $j < i$
- This notion is generalized to k-records: $\sigma(i)$ is a k-record if there are at most $k-1$ elements $\sigma(j)$ larger than $\sigma(i)$ for $j < i$; in other words, $\sigma(i)$ is among the $k$ largest elements in $\sigma(1), \ldots, \sigma(i)$

---
Example

This example permutation contains six 2-records

$$\mathcal{P} = 3, 6, 1, 12, 8, 10, 4, 13, 7, 5, 9, 11, 2$$

# Recordinality

- $\sigma(i)$ is a record of the permutation $\sigma$ if $\sigma(i) > \sigma(j)$ for all $j < i$
- This notion is generalized to k-records: $\sigma(i)$ is a k-record if there are at most $k - 1$ elements $\sigma(j)$ larger than $\sigma(i)$ for $j < i$; in other words, $\sigma(i)$ is among the k largest elements in $\sigma(1), \ldots, \sigma(i)$

---
Example

This example permutation contains six 2-records

$$\mathcal{P} = 3, 6, 1, 12, 8, 10, 4, 13, 7, 5, 9, 11, 2$$
---

# Recordinality

```
procedure Recordinality(𝒵, k)
    fill 𝒮 with the first k distinct elements (hash values)
    of the stream 𝒵
    R ← k
    for all z ∈ 𝒵 do
        y ← h(z)
        if y > min{h(x) | x ∈ 𝒮} ∧ z ∉ 𝒮 then
            z* ← the element in 𝒮 with min. hash value
            R ← R + 1; 𝒮 ← 𝒮 ∪ {z} \ z*
        end if
    end for
    return  Z = k (1 + 1/k)^(R−k+1) − 1
end procedure
```

Memory: $k$ hash values ($k \log n$ bits) + 1 counter ($\log \log n$ bits)

## Analysis of k-Records

The behavior of $R = R_n$, the number of $k$-records in a random permutation of size $n$, is very well understood[1]

$$E[R] = k(H_n - H_k + 1) = k \ln(n/k) + O(1)$$

Likewise

$$\text{Var}[R] = k(H_n - H_k) - k^2(H_n^{(2)} - H_k^{(2)}) = k \ln(n/k) + O(1)$$

and we also know exact and asymptotic estimates for $\text{Prob}\{R = j\}$.

---

[1] $H_n = 1 + 1/2 + 1/3 + \cdots + 1/n \sim \ln n + O(1)$ denotes the $n$-th harmonic number, and $H_n^{(2)} = 1 + 1/4 + 1/9 + \cdots + 1/n^2 \leqslant \pi^2/6$.

# The Estimator for Recordinality

Let us assume for the moment that $k \leqslant R \leqslant n$. If $R < k$ then we are sure that $n = R$. Otherwise, since $E[R] = k \ln(n/k) + O(1)$ we can take

$$Z = \exp(\phi \cdot R)$$

for some correcting factor $\phi$ to be determined and such that $E[Z]$ is (asymptotically?) $n$. Our knowledge of the probability distribution of $R$ furnishes the exact form for $Z$.

# The Estimator for Recordinality

> **Theorem**
>
> Let $R$ be the number of $k$-records seen while processing the data stream $\mathcal{Z}$. Then
>
> $$Z := k \left( 1 + \frac{1}{k} \right)^{R-k+1} - 1$$
>
> is an unbiased estimator of the cardinality (number of distinct elements) of $\mathcal{Z}$, that is,
>
> $$E[Z] = n$$

# Part IV

# Distinct Sampling and Applications

# Drawing Random Samples



- In a random sample from the data stream (e.g., using the reservoir method) each distinct element $x_j$ appears with relative frequency in the sample equal to its relative frequency $f_j/N$ in the data stream $\Rightarrow$ needle-on-a-haystack

- Elements of low frequency will seldom be sampled, and we cannot keep exact counts as we don't know if the sampled elements have been "monitorized" from the beginning

# Drawing Random Samples



- In a random sample from the data stream (e.g., using the reservoir method) each distinct element $x_j$ appears with relative frequency in the sample equal to its relative frequency $f_j/N$ in the data stream $\Rightarrow$ needle-on-a-haystack
- Elements of low frequency will seldom be sampled, and we cannot keep exact counts as we don't know if the sampled elements have been "monitorized" from the beginning

# Drawing Random Samples



- The distinct sampling problem is to draw a random sample of distinct elements and it has many applications in data stream analysis

- For example, to estimate the number of k-elephants or k-mice in the stream we can draw a random sample of S distinct elements, together with their frequency counts

- Let $S_P$ be the number of mice (or elephants) in the sample, and $n_P$ the number of mice (or elephants) in the data stream. Then

$$E\left[\frac{S_P}{S}\right] = \frac{n_P}{n}$$

# Drawing Random Samples



- The distinct sampling problem is to draw a random sample of distinct elements and it has many applications in data stream analysis

- For example, to estimate the number of k-elephants or k-mice in the stream we can draw a random sample of $S$ distinct elements, together with their frequency counts

- Let $S_P$ be the number of mice (or elephants) in the sample, and $n_P$ the number of mice (or elephants) in the data stream. Then

$$E\left[\frac{S_P}{S}\right] = \frac{n_P}{n}$$

# Drawing Random Samples

- The distinct sampling problem is to draw a random sample of distinct elements and it has many applications in data stream analysis

- For example, to estimate the number of k-elephants or k-mice in the stream we can draw a random sample of $S$ distinct elements, together with their frequency counts

- Let $S_P$ be the number of mice (or elephants) in the sample, and $n_P$ the number of mice (or elephants) in the data stream. Then

$$E\left[\frac{S_P}{S}\right] = \frac{n_P}{n}$$

# Drawing Random Samples

Let P some property.

- $n = \#$ of distinct elements in $\mathcal{Z}$
- $n_P = \#$ of distinct elements in $\mathcal{Z}$ that satisfy P
- $S =$ size of the sample $\Leftarrow$ in general, a r.v., assume $2 \leqslant S \leqslant n$
- $S_P = \#$ of elements in the sample that satisfy P

**Theorem**

**1** $E\left[\dfrac{S_P}{S}\right] = \dfrac{n_P}{n}$

**2** $Var\left[\dfrac{S_P}{S}\right] \sim \dfrac{n_P}{n} \cdot \left(1 - \dfrac{n_P}{n}\right) \cdot E\left[\dfrac{1}{S}\right]$

# Drawing Random Samples

Let P some property.

- $n = \#$ of distinct elements in $\mathcal{Z}$
- $n_P = \#$ of distinct elements in $\mathcal{Z}$ that satisfy P
- $S = $ size of the sample $\Leftarrow$ in general, a r.v., assume $2 \leqslant S \leqslant n$
- $S_P = \#$ of elements in the sample that satisfy P

---
**Theorem**

**1** $E\left[\dfrac{S_P}{S}\right] = \dfrac{n_P}{n}$

**2** $Var\left[\dfrac{S_P}{S}\right] \sim \dfrac{n_P}{n} \cdot \left(1 - \dfrac{n_P}{n}\right) \cdot E\left[\dfrac{1}{S}\right]$

---

# Adaptive Sampling



M. Wegman    G. Louchard

- *Adaptive sampling* (Wegman, 1980; Flajolet, 1990; Louchard *et al*, 1997) is the first algorithm proposed specifically for distinct sampling

- It also gives an estimation of the cardinality, as the size $S$ of the returned sample is itself a random variable, but it is always bounded by a fixed constant $maxS$

## Adaptive Sampling

```
procedure AdaptiveSampling(𝒵, maxS)
    𝒮 ← ∅; p ← 0
    for z ∈ 𝒵 do
        if hash(z) = 0^p . . . ∧ z ∉ 𝒮  then
            𝒮 ← 𝒮 ∪ {z}
            if |𝒮| > maxS then
                p ← p + 1
                𝒮 ← 𝒮 \ {z ∈ 𝒮 | h(z) = 0^{p−1}1 . . .} // Filter 𝒮
            end if
        end if
    end for
    return 𝒮
end procedure
```

The set $\mathcal{S}$ is a random sample (because we can assume hash values behave as random uniform numbers) of $S = |\mathcal{S}|$ distinct elements; if $n$ is large enough, $maxS/2 \leqslant \mathrm{E}\,[S] \leqslant maxS$

# Adaptive Sampling

At the end of the algorithm, $S$ is the number of distinct elemnts with hash value starting $.0^p \equiv$ the number of strings in the subtree rooted at $0^p$ in a binary trie for $n$ random binary strings. There are $2^p$ subtrees rooted at depth $p$

$$S = |\mathcal{S}| \approx n/2^p \Rightarrow E\left[2^p \cdot S\right] \approx n$$

# Distinct Sampling in Recordinality and Order Statistics

- Recordinality and KMV collect the elements with the $k$ largest (smallest) hash values

- Such $k$ elements constitute a random sample of $k$ distinct elements, because hash values behave as random numbers; but the value $k$ is fixed in advance and might be too small for the sample to be representative

- Recordinality can be easily adapted to collect random samples of expected size $\Theta(\log n)$ or $\Theta(n^{\alpha})$, with $0 < \alpha < 1$ and without prior knowledge of $n$! $\Rightarrow$ Affirmative Sampling $\Rightarrow$ variable-size samples, growing with $n$, better precision in inferences about the full data stream

# Distinct Sampling in Recordinality and Order Statistics

- Recordinality and KMV collect the elements with the $k$ largest (smallest) hash values
- Such $k$ elements constitute a random sample of $k$ distinct elements, because hash values behave as random numbers; but the value $k$ is fixed in advance and might be too small for the sample to be representative
- Recordinality can be easily adapted to collect random samples of expected size $\Theta(\log n)$ or $\Theta(n^\alpha)$, with $0 < \alpha < 1$ and without prior knowledge of $n$! $\Rightarrow$ Affirmative Sampling $\Rightarrow$ variable-size samples, growing with $n$, better precision in inferences about the full data stream

# Distinct Sampling in Recordinality and Order Statistics

- Recordinality and KMV collect the elements with the $k$ largest (smallest) hash values
- Such $k$ elements constitute a random sample of $k$ distinct elements, because hash values behave as random numbers; but the value $k$ is fixed in advance and might be too small for the sample to be representative
- Recordinality can be easily adapted to collect random samples of expected size $\Theta(\log n)$ or $\Theta(n^\alpha)$, with $0 < \alpha < 1$ and without prior knowledge of $n$! $\Rightarrow$ Affirmative Sampling $\Rightarrow$ variable-size samples, growing with $n$, better precision in inferences about the full data stream

# Affirmative Sampling



- Early ideas date back to the original paper on Recordinality (2012); developed and analyzed in detail in (Lumbroso, M., 2019, 2022)
- The larger the cardinality ($n$) the larger the samples $\Rightarrow$ samples better represent the population
- All distinct elements have the same opportunity to be sampled, and if sampled they can be "monitorized" from their first appearance

# Affirmative Sampling

```
procedure AffirmativeSampling(k, 𝒵)
    fill 𝒮 with the first k distinct elements (and hash values)
      of the stream 𝒵
    for all y ∈ 𝒵 do
        z ← hash(y)
        if z < z* then // z* = min hash in 𝒮 = hash(y*)
            Discard y
        else if y ∈ 𝒮 then
            Update y stats
        else if z > z^(k) then // z^(k) = k-th largest hash in 𝒮
            S ← S ∪ {y} // Add y to the sample
        else
            S ← S \ {y*} ∪ {y} // Replace y* by y in the sample
        end if
    end for
    return 𝒮
end procedure
```

# Affirmative Sampling



$z > z^{(k)} \rightarrow$ add $z$ to $\mathcal{S}$ ($\mathcal{S}$ grows by 1)

$z^{(k)} = k$-th largest hash in $\mathcal{S}$
$\quad = k$-th largest hash in the data stream $\mathcal{Z}$

$z^{(k)} > z > z^* \rightarrow$ replace $z^*$ by $z$ in $\mathcal{S}$

$z^* = $ minimum hash in $\mathcal{S}$
$\quad = |S|$-th largest hash in the data stream $\mathcal{Z}$

$z \in \mathcal{S} \vee z < z^* \rightarrow$ update/discard $z$

$S = $ number of left-to-right $k$-records
$\quad = $ size of $\mathcal{S}$
$E[S] \sim k \ln(n/k) + k$

# Affirmative Sampling

- The size $S$ of the sample $\mathcal{S}$ is a random variable $=$ the number of $k$-records in a random permutation of size $n \Rightarrow$ $E[S] = k \ln(n/k) + \mathcal{O}(1)$

- The sample does not contain the $k$-records, but the $S$ elements with the largest hash values seen so far $\Rightarrow$ $\mathcal{S}$ is a random sample

- If $x \in \mathcal{S}$ then $x$ has been added to $S$ in its very first occurrence and it has remained in $\mathcal{S}$ ever since $\Rightarrow$ can collect exact stats (e.g. frequency counts) for $x$

# Affirmative Sampling

- We also understand fairly well $F$ = number of times an element substitutes another in the sample (not a $k$-record, but larger than some $k$-record):

$$E[F] = k \ln^2(n/k) + \text{l.o.t.}$$

- Expected cost $C_{N,n}$ of Affirmative Sampling

$$E[C_{N,n}] = \Theta\left(N + k \log^2(n/k) \log\log(n/k)\right)$$

using appropriate data structures for the sample $\mathcal{S}$

# Similarity Estimation

Consider two data streams $\mathcal{Z}_A$ and $\mathcal{Z}_B$. Let $A$ and $B$ denote their respective sets of distinct elements. Similarity between the two sets is often measured by their Jaccard index

$$J(A, B) = \frac{|A \cap B|}{|A \cup B|}$$

The containment index measures how much "$A \subseteq B$" and it is given by

$$c(A, B) = \frac{|A \cap B|}{|A|}$$

# Similarity Estimation

We can estimate similarity and containment from random samples $S_A$ and $S_B$ of the two streams. If the samples are drawn using Affirmative Sampling then

---
**Theorem**

1. $E\left[J(S'_A, S'_B)\right] = J(A, B) = \frac{|A \cap B|}{|A \cup B|}$

2. $Var\left[J(S'_A, S'_B)\right] \sim \frac{J(A,B) \cdot (1 - J(A,B))}{k \ln(|A \cup B|/k)}$
---

# Similarity Estimation



$A$

$S_A'$
$= S_A$

$S_B'$

$S_B$

$B$

$S_A' \cap S_B' = S_A \cap S_B$

$= (S_A' \cup S_B') \cap (A \cap B)$

$S_A' \cup S_B'$ is
a random sample
of $A \cup B$ ($S_A \cup S_B$
is not!)

# Estimating the size of the intersection

We can estimate the size of the intersection with:

$$Z_1 = \frac{|S_A \cap S_B|}{|S_A|} \cdot \left( k \left( 1 + \frac{1}{k} \right)^{|S_A| - k + 1} - 1 \right)$$

$$Z_2 = \frac{|S_A \cap S_B|}{|S_A|} \cdot \frac{|S_A| - 1}{1 - M_{S_A}}, \qquad M_{S_A} = \min\{h(z) \mid z \in S_A\}$$

$$E[Z_1] = E[Z_2] = |A \cap B|$$

N.B. No need to "filter" the samples

# Other similarity measures

| Jaccard's index | $\frac{|A \cap B|}{|A \cup B|}$ |
|---|---|
| Otsuka-Ochiai (a.k.a. Cosine) | $\frac{|A \cap B|}{\sqrt{|A| \cdot |B|}}$ |
| Sørensen-Dice | $2\frac{|A \cap B|}{|A| + |B|}$ |
| Kulczynski 1 | $\frac{|A \cap B|}{|A \triangle B|}$ |
| Kulczynski 2 | $\frac{1}{2}\left(\frac{|A \cap B|}{|A|} + \frac{|A \cap B|}{|B|}\right)$ |
| Simpson | $\frac{|A \cap B|}{\min(|A|, |B|)}$ |
| Braun-Blanquet | $\frac{|A \cap B|}{\max(|A|, |B|)}$ |
| Correlation | $\cos^2(A, B) = \frac{|A \cap B|^2}{|A| \cdot |B|}$ |
| ... | ... |

# Other similarity measures

The same proof that works for Jaccard's similarity also works for containment and many other similarity measures:

**1** $\mathsf{E}\left[c(S_A, S_B)\right] = c(A, B) = |A \cap B|/|A|$

**2** If $\sigma$ is any of Jaccard, Simpson, Braun-Blanquet, Kulczynski 2, correlation or Sørensen-Dice:

$$\mathsf{E}\left[\sigma(S'_A, S'_B)\right] = \sigma(A, B)$$

**3** It also works for cosine and Kulczynski 1 similarities and many others because these distances can be expressed as $f(J(A, B))$; while $\mathsf{E}\left[f(X)\right] \neq f(\mathsf{E}\left[X\right])$ one can show that $\mathsf{E}\left[f(X)\right] \sim f(\mathsf{E}\left[X\right])$ when we use samples of variable size to estimate $J(A, B)$, since the variance and all central moments of the estimator $\to 0$ as $\min(A, B) \to \infty$

# Conclusions

- Neeeded: easy and practical algorithms, often randomized, precise mathematical analysis is a must
- We have the right arsenal of tools, there is plenty of open problems in data streaming for which we might have a say
- Many elegant and challenging mathematical problems
- Real-life applications, right motivations and incentives $\implies$ practical relevant algorithms used by thousands of practitioners on a daily basis (e.g., HyperLogLog is part of the "infrastructure" of all major data analytics companies)
- Not by chance:
    - Flajolet pioneered some of the most important techniques and results (and his most cited works are those he did in this area)
    - Two Flajolet Lecturer awardees, Sedgewick and Janson, join forces for HyperBit, the ultimate(?) cardinality estimator

# To Know More

[1] Ziv Bar-Yossef, T. S. Jayram, Ravi Kumar, D. Sivakumar, and Luca Trevisan.
Counting Distinct Elements in a Data Stream.
*Randomization and Approximation Techniques (RANDOM)*, pages 1–10. 2002.

[2] Andrei Broder.
On the resemblance and containment of documents.
*Proc. Compression and Complexity of Sequences (SEQUENCES)*, pages 21–29. 1997.

[3] Marianne Durand and Philippe Flajolet.
LogLog Counting of Large Cardinalities.
*Proc. European Symposium on Algorithms (ESA)*, volume 2832 of *Lecture Notes in Computer Science*, pages 605–617, 2003.

# To Know More

[4]   Philippe Chassaing and Lucas Gerin.
      Efficient Estimation of the Cardinality of Large Data Sets.
      *Proc. Int. Col. Mathematics and Computer Science
      (MathInfo)*, pages 419–422, 2007.

[5]   Philippe Flajolet.
      On adaptive sampling.
      *Computing*, 34:391–400, 1990.

[6]   Philippe Flajolet, Éric Fusy, Olivier Gandouet, and Frédéric
      Meunier.
      HyperLoglog: the analysis of a near-optimal cardinality
      estimation algorithm.
      *Proceedings of Int. Conf. Analysis of Algorithms (AofA)*, pages
      127–146, 2007.

# To Know More

[7]   Philippe Flajolet and G. Nigel Martin.
      Probabilistic Counting Algorithms for Data Base Applications.
      *Journal of Computer and System Sciences*, 31(2):182–209,
      1985.

[8]   A. Helmi, J. Lumbroso, C. Martínez, and A. Viola.
      Counting distinct elements in data streams: the random
      permutation viewpoint.
      *Proc. of Int. Conf. Analysis of Algorithms (AofA)*, pages
      323–338, 2012.

[9]   Jérémie Lumbroso.
      An optimal cardinality estimation algorithm based on order
      statistics and its full analysis.
      In *Proc. Analysis of Algorithms (AofA)*, pages 489–504, 2010.

## To Know More

[10] Jérémie Lumbroso.
How Flajolet Processed Stream with Coin Flips.
arXiv:1805.00612v1 [cs.DS] 2 May 2018. December 2013.

[11] Jérémie Lumbroso and Conrado Martínez.
Affirmative Sampling: Theory and Applications.
In *Proc. 33rd Analysis of Algorithms (AofA)*, LIPIcs vol. 225, pages 12:1–12:17, 2022.

# To Know More

[12] M. Monenizadeh and D. Woodruff.
*1-Pass Relative-Error $L_p$-Sampling with Applications*.
In *Proc. Symp. Discrete Algorithms (SODA)*, pages
1143–1160, 2010.

[13] S. Muthu Muthukrishnan.
Data streams: Algorithms and applications.
*Foundations and Trends in Theoretical Computer Science*,
1(2):117–236, 2005.