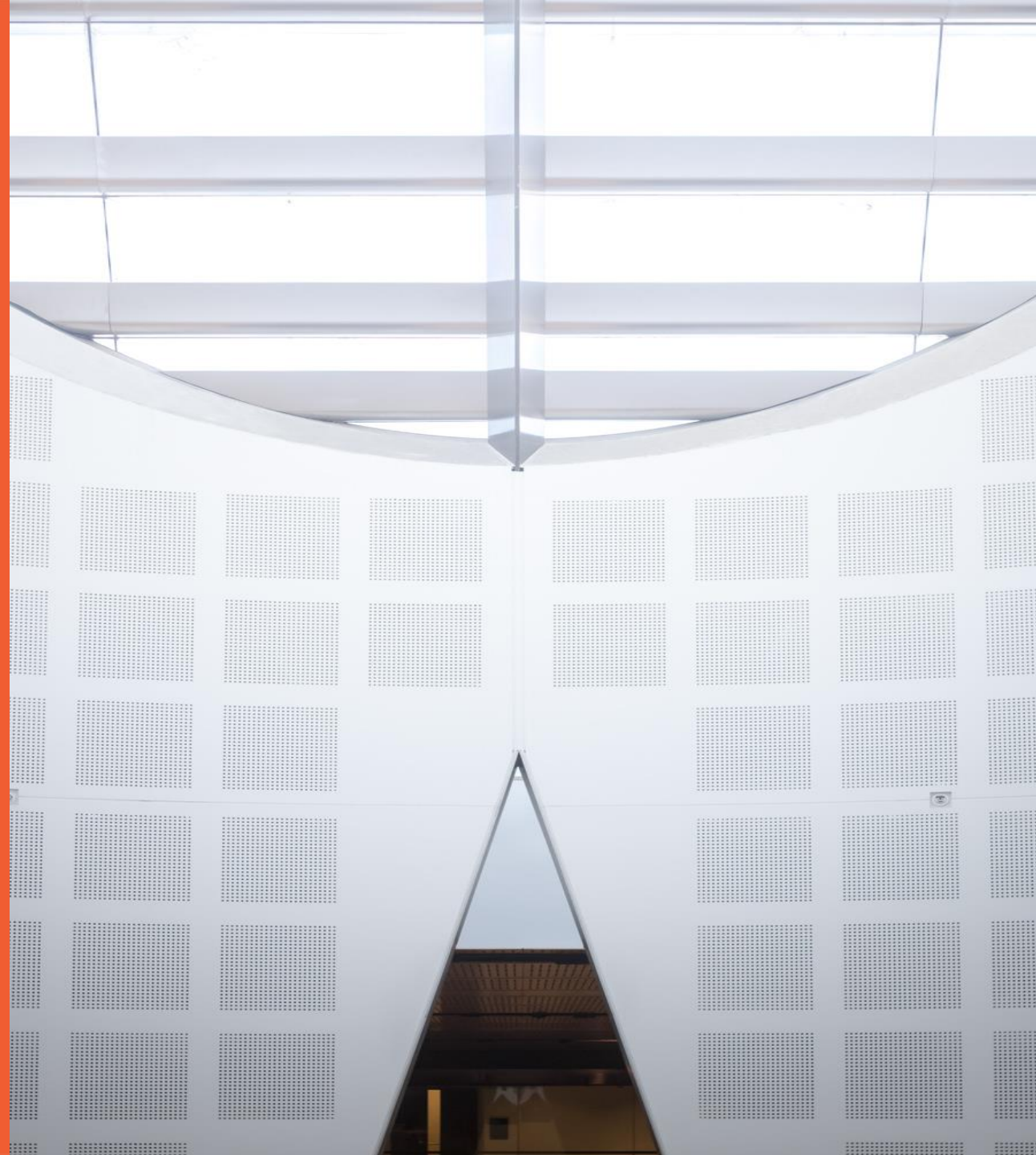COMPx270: Randomised and Advanced Algorithms
Lecture 8:  Streaming and Sketching I

Clément Canonne

School of Computer Science

THE UNIVERSITY OF
SYDNEY

# Some housekeeping

- A2 due tonight

    See Ed+email announcement about Q3.f

- A3 now live, due May 9

- No class next week (semester break!)

# A question

You have a graph, coming one edge at a time, with possible duplicates, and no paper to write anything done, only your memory. What is its average degree?

# A question

You have a graph, coming one edge at a time, with possible duplicates, and no paper to write anything done, only your memory. What is its average degree?

(1,2)

# A question

You have a graph, coming one edge at a time, with possible duplicates, and no paper to write anything done, only your memory. What is its average degree?

(2,4)

# A question

You have a graph, coming one edge at a time, with possible duplicates, and no paper to write anything done, only your memory. What is its average degree?

(1,2)

# A question

You have a graph, coming one edge at a time, with possible duplicates, and no paper to write anything done, only your memory. What is its average degree?

(4,5)

# A question

You have a graph, coming one edge at a time, with possible duplicates, and no paper to write anything done, <span style="color:blue">only your memory</span>. <span style="color:red">What is its average degree?</span>

<p style="text-align:center;color:red">(4,5)</p>

# A question

You have a graph, coming one edge at a time, with possible duplicates, and no paper to write anything done, <span style="color:blue">only your memory</span>. <span style="color:darkred">What is its average degree?</span>

<span style="color:darkred">(3,4)</span>

# A question

You have a graph, coming one edge at a time, with possible duplicates, and no paper to write anything done, <span style="color:blue">only your memory</span>. <span style="color:red">What is its average degree?</span>

<span style="color:red">(3,6)</span>

# A question

You have a graph, coming one edge at a time, with possible duplicates, and no paper to write anything done, only your memory. What is its average degree?

(1,4)

# A question

You have a graph, coming one edge at a time, with possible duplicates, and no paper to write anything done, only your memory. What is its average degree?

(4,6)

# A question

You have a graph, coming one edge at a time, with possible duplicates, and no paper to write anything done, only your memory. What is its average degree?

(3,5)

# A question

You have a graph, coming one edge at a time, with possible duplicates, and no paper to write anything done, only your memory. What is its average degree?

(3,4)

# A question

You have a graph, coming one edge at a time, with possible duplicates, and no paper to write anything done, only your memory. What is its average degree?

(4,5)

# A question

You have a graph, coming one edge at a time, with possible duplicates, and no paper to write anything done, only your memory. What is its average degree?

(3,4)

# A question

You have a graph, coming one edge at a time, with possible duplicates, and no paper to write anything done, only your memory. What is its average degree?

(3,6)

# A question

You have a graph, coming one edge at a time, with possible duplicates, and no paper to write anything done, only your memory. What is its average degree?

(1,4)

# A question

You have a graph, coming one edge at a time, with possible duplicates, and no paper to write anything done, only your memory. What is its average degree?
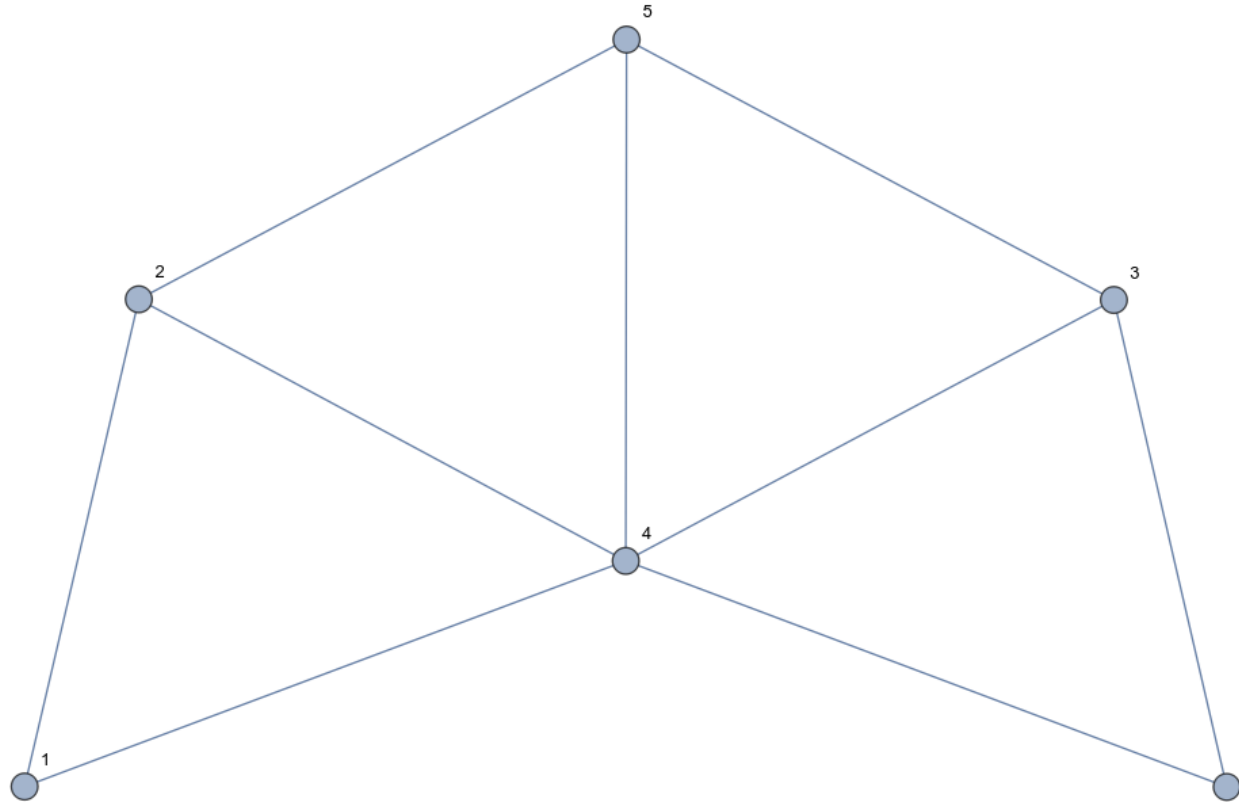
(4,6)

# A question (an answer)

# Streaming algorithms: what? (1/3)

. Low memory : cannot store whole input $\nearrow$ so less than $O(m \log n)$ or $O(n \log m)$

. Input comes as a <u>stream</u> : sequence of length $\textcircled{m}$

$$\sigma = (a_1, a_2, -, a_m)$$
$$\uparrow$$
$$a_i \in \mathcal{X} , |\mathcal{X}| = \textcircled{n}$$

Worst-case (arbitrary) order.


$$\underset{p}{\overset{\longrightarrow}{\vdash \cdots \cdots \dashv}} \sigma$$

. $p$-pass algorithms get to see $\sigma$ $p$ times $(p=1$ for us$^*)$

. <u>cash register</u> : don't remove parts of the input (that would be "turnstile")

SPACE : $o(\min(n,m))$
("sublinear")

hope : $O(\log(mn)) = O(\log m + \log n)$
very good : $polylog(m,n)$

# Streaming algorithms: what? (3/3)

- Randomised algorithms
- Approximate: want to compute some value $v \geq 0$

we're OK with $\hat{v} \overset{\varepsilon}{\approx} v$

$\uparrow$ ?

① Multiplicative:
$$\Pr[\, |\hat{v} - v| \geq \varepsilon v \,] \leq \delta$$

$(1 \pm \varepsilon)$ approx

② Additive:
$$\Pr[\, |\hat{v} - v| \geq \varepsilon \,] \leq \delta$$

$\pm \varepsilon$ approx

# First example: Majority

A.k.a. "special case of Heavy Hitters"

- MAJORITY : is there an elem$^t$ appearing $\geq 50\%$ of the time in the stream? (If so, which one(s)?)

$$\sigma = (\sigma_1, \ldots, \sigma_m) \in [n]^m$$

$$\forall j \in [n] \qquad f_j = \#\left(\begin{array}{c}\text{times } j \text{ appears} \\ \text{in } \sigma\end{array}\right) = \sum_{i=1}^{m} \mathbb{1}_{\sigma_i = j}$$

<span style="color:red">frequency of element $j$</span>

<span style="color:red">Frequency vector</span> $\vec{f} = (f_1, \ldots, f_n)$

- $0 \leq f_j \leq m \quad \forall j \in [n]$

- $F_1 = \|\vec{f}\|_1 = \sum_{j=1}^{n} f_j = m$

# First example: Majority (Frequency Estimation)

MAJORITY:  "Is there $j \in [n]$ st. $f_j \geq \frac{m}{2}$?"  (at most 2 of them)

$\varepsilon$-HH:  "Is there $j \in [n]$ s.t. $f_j \geq \varepsilon m$?"  (at most $\frac{1}{\varepsilon}$ of them)

Want to solve this in one pass.

We'll see two passes, but <u>deterministically</u>

# First example: the Misra-Gries algorithm (1/3)

MISRA-GRIES    returns $\hat{b}_1, \dots, \hat{b}_n$    (a succinct representation of them)

st.    $b_j - \varepsilon m \leq \hat{b}_j \leq b_j \quad \forall j$

in <u>one</u> pass.

⚠ only $O(\frac{1}{\varepsilon})$ estimates are non-zero: only returns these $\hat{b}_j$

→ TO SOLVE MAJORITY IN TWO PASSES

Pass ① Run M-G on $\sigma$ with $\varepsilon = \frac{1}{4}$

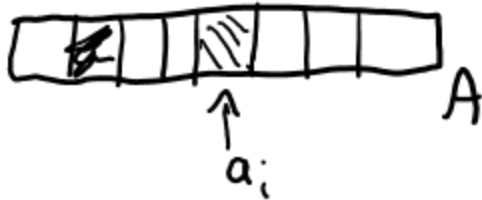Pass ② Count exactly the frequency for all $j$'s s.t. $\hat{b}_j \geq \frac{m}{4}$

If $b_j \geq \frac{m}{2}$

$\hat{b}_j \geq \frac{m}{4} > 0$

If $\hat{b}_j \geq \frac{m}{4}$

$b_j \geq \frac{m}{4}$

→ at most 4 of them

# First example: the Misra-Gries algorithm, alternative view (2/3)

$A \leftarrow n$ zeroes       (use a BST to save space)

$k \leftarrow 1/\varepsilon$

At step $1 \le i \le m$:

- get $a_i$

- If $A[a_i] > 0$

  $A[a_i] += 1$        # of non-zero entries.

- If $A[a_i] = 0$ and $|A| < k-1$

  $A[a_i] = 1$

- If $A[a_i] = 0$ and $|A| = k-1$

  $A[a_i]=1 \nearrow$ For all $j$ s.t $A[j] > 0$

  $A[j] = A[j] - 1$



A

At the end : return all $j$'s (and $A[j]$) s.t. $A[j] > 0$

SPACE : $O(k \cdot (\log m + \log n))$

$= O\left(\frac{\log(nm)}{\varepsilon}\right)$

CORRECTNESS     $(\hat{f}_j = A[j])$

① $\forall j, \quad \hat{f}_j \le f_j$      "clear"

② $\forall j \quad \hat{f}_j \ge f_j - \frac{m}{k}$

Claim: "can't decrement an element too many times"

Each decrement corresponds to exactly $k$ prior increments $\rightarrow$ at most $\frac{m}{k}$ decrements

# First example: the Misra-Gries algorithm (3/3)

**Theorem 39.** *The* MISRA-GRIES *algorithm is a deterministic one-pass algorithm which, for any given parameter $\varepsilon \in (0,1]$, provides $\hat{f}_1, \ldots, \hat{f}_n$ of all element frequencies such that*

$$f_j - \varepsilon m \leq \hat{f}_j \leq f_j, \qquad j \in [n]$$

*with space complexity $s = O(\log(mn)/\varepsilon)$. (In particular, it can be used to solve the* MAJORITY *problem in two passes.)*

# Second example: Approximate Counting

$$n = 2 \qquad \mathcal{X} = \{0, 1\} \qquad \text{Want} \quad d = \sum_{i=1}^{m} a_i$$

- $O(\log m)$ trivial : counter (deterministic)
- 2-estimate $O(\log \log m)$ space (Morris)

randomized

# Second example: Approximate Counting and the Morris Counter

1: $x \leftarrow 0$

2: **for all** $1 \leq i \leq m$ **do**

3:       Get item $a_i \in \{0,1\}$

4:       **if** $a_i = 1$ **then**

5:            $r_i \leftarrow \text{Bern}(1/2^x)$         ▷ Independent of previous choices.

6:            $x \leftarrow x + r_i$

7: **return** $\widehat{d} \leftarrow 2^x - 1$

$$E[\widehat{d}] = ?$$
$$\text{Var}[\widehat{d}] = ?$$

# Second example: Approximate Counting and the Morris Counter

1: $C_0 \leftarrow 1$

2: **for all** $1 \le i \le m$ **do**

3:      Get item $a_i \in \{0,1\}$

4:      **if** $a_i = 1$ **then**

5:          $r_i \leftarrow \mathrm{Bern}(1/C_{i-1})$      $\triangleright$ Independent of previous choices.

6:      **else** $r_i \leftarrow 0$

7:      $C_i \leftarrow 2^{r_i} C_{i-1}$

8: **return** $\hat{d} \leftarrow C_m - 1$

$$\text{Claim}: \qquad \mathbb{E}[C_m] = d+1$$

$$\mathrm{Var}[C_m] = O(d^2) \qquad \left(= \binom{d}{2}\right)$$

# Throwback: Law of Total Expectation (and Friends)

$$\mathbb{E}[\mathbb{E}[X|Y]] = \mathbb{E}[X]$$

$$\mathbb{E}[\beta(Y)|Y] = \beta(Y)$$

$$\mathbb{E}[\beta(Y)X|Y] = \beta(Y)\,\mathbb{E}[X|Y]$$

# Second example: the Morris Counter (1/3)

1:   $C_0 \leftarrow 1$
2:   **for all** $1 \le i \le m$ **do**
3:     Get item $a_i \in \{0,1\}$
4:     **if** $a_i = 1$ **then**
5:       $r_i \leftarrow \mathrm{Bern}(1/C_{i-1})$    ▷ Independent of previous choices.
6:     **else** $r_i \leftarrow 0$
7:     $C_i \leftarrow 2^{r_i} C_{i-1}$
8: **return** $\hat{d} \leftarrow C_m - 1$

- $\mathbb{E}[C_0] = C_0 = 1$

- $\mathbb{E}[C_{i+1} \mid C_i] \stackrel{?}{=} \quad$ $\left( \text{If } a_{i+1} = 0, \quad C_{i+1} = C_i \right)$

$$\frac{1}{C_i} \cdot 2C_i + \left(1 - \frac{1}{C_i}\right) \cdot C_i$$

$\uparrow$
$\text{If } a_{i+1} = 1$

$$= C_i + 1$$

$$\mathbb{E}[C_{i+1} \mid C_i] = C_i + a_{i+1}$$

$\rightarrow \quad \mathbb{E}[C_{i+1}] = \mathbb{E}[\mathbb{E}[C_{i+1} \mid C_i]] = \mathbb{E}[C_i] + a_{i+1} = \mathbb{E}[C_{i-1}] + a_i + a_{i+1} = \dots$

$\rightarrow \quad \mathbb{E}[C_m] = 1 + \sum_{i=1}^{m} a_i = 1 + d$  ✓

# Second example: the Morris Counter (2/3)

1: $C_0 \leftarrow 1$
2: **for all** $1 \leq i \leq m$ **do**
3:     Get item $a_i \in \{0,1\}$
4:     **if** $a_i = 1$ **then**
5:        $r_i \leftarrow \text{Bern}(1/C_{i-1})$     $\triangleright$ Independent of previous choices.
6:     **else** $r_i \leftarrow 0$
7:     $C_i \leftarrow 2^{r_i} C_{i-1}$
8: **return** $\hat{d} \leftarrow C_m - 1$

$$\text{Var}[C_m] = \mathbb{E}[C_m^2] - \mathbb{E}[C_m]^2$$

$$\underset{\text{known!}\ (d+1)^2}{\uparrow}$$

$$\mathbb{E}[C_m^2] = ?$$

$$\cdot \ \mathbb{E}[C_0^2] = C_0^2 = 1$$

$$\cdot \ \mathbb{E}[C_{i+1}^2 \mid C_i] = \begin{cases} C_i^2 & \text{if } a_{i+1} = 0 \\ \dfrac{1}{C_i} \cdot 4C_i^2 + \left(1 - \dfrac{1}{C_i}\right) C_i^2 = 3C_i + C_i^2 & \text{if } a_{i+1} = 1 \end{cases}$$

$$= C_i^2 + a_{i+1}(2 + a_{i+1}) C_i$$

$$+ \sum_{i=1}^{m} (\ldots)$$

$$\rightarrow \quad \mathbb{E}[C_m^2] = \mathbb{E}[\mathbb{E}[C_m^2 \mid C_{m-1}]] = \cdots = 1 + 3\frac{d(d+1)}{2}$$

$$\rightsquigarrow \quad \text{Var } C_m = \frac{d(d-1)}{2} \quad \checkmark$$

$$\mathbb{E}[C_m] = d+1 \qquad \checkmark$$

$$\text{Var}[C_m] = \Theta(d^2) \qquad \times$$

Chebyshev:

$$C_m = \mathbb{E}[C_m] \pm \sqrt{\text{Var}[C_m]}$$

w.p. 99%

useless guarantee.

Doom?   No

① "Meh" guarantee,  ✓
     w.p. 51%

② "Good" guarantee,
     take average  w.p. 51%

③ "Good" guarantee
     w.p. 1−δ
     median trick

# Second example: the Morris Counter, Median-of-Means

**Theorem 40.** *The medians-of-means version of the* MORRIS COUNTER *is a* randomised *one-pass algorithm which, for any given parameters* $\varepsilon, \delta \in (0,1]$, *provides an estimate* $\widehat{d}$ *of the number* $d$ *of non-zero elements of the stream such that*

$$\Pr\left[ (1-\varepsilon)d \leq \widehat{d} \leq (1+\varepsilon)d \right] \geq 1 - \delta$$

*with space complexity*

$$s = O\left( \frac{\log\log m}{\varepsilon^2} \cdot \log\frac{1}{\delta} \right)$$

*that is,* doubly logarithmic *in* $m$.

③ Median truck over
$T' = O(\log \frac{1}{\delta})$ instances

average over

② $T = O(\frac{1}{\varepsilon^2})$ instances to reduce variance , then Chebyshev

# Did we need to do that?

No.

No need for median-of-means here!

⊕ better space ...

## Second example: the Morris Counter, careful version (1/2)

Morris from before

$$C \leftarrow 2C \text{ wp } \frac{1}{C} \quad (\text{when } a_i = 1)$$

Morris, better

$$C \leftarrow (1+\alpha)C \text{ wp } \frac{1}{\alpha C}$$

$$\boxed{\alpha = 2\varepsilon^2 \delta}$$

$$(\text{Estimate: } (1+\alpha)^x - 1)$$

# Second example: the Morris Counter, careful version (2/2)

**Theorem 41.** *The "careful" version of* Morris Counter *is a randomised one-pass algorithm which, for any given parameters* $\varepsilon, \delta \in (0,1]$, *provides an estimate* $\widehat{d}$ *of the number* $d$ *of non-zero elements of the stream such that*

$$\Pr\left[ (1 - \varepsilon)d \leq \widehat{d} \leq (1 + \varepsilon)d \right] \geq 1 - \delta$$

*with space complexity*

$$s = O\left( \log\log m + \log\frac{1}{\varepsilon} + \log\frac{1}{\delta} \right)$$

*that is, doubly logarithmic in* $m$ *and logarithmic in* $1/\varepsilon$.

# Third example: Distinct Elements

$$\text{Approximate } F_0 = \underbrace{\sum_{j=1}^{n} \mathbb{1}_{f_i > 0}}_{\text{"d" (new name)}}$$

$$\text{Return } \hat{d} = (1 \pm \varepsilon) d$$

$$\varepsilon = \Theta(1)$$

---

1: Pick $h \colon [n] \to [n]$ from a strongly universal hashing family
2: $z \leftarrow 0$
3: **for all** $1 \le i \le m$ **do**
4:      Get item $a_i \in [n]$
5:      **if** zeros($h(a_i)$) $\ge z$ **then**
6:          $z \leftarrow$ zeros($h(a_i)$)
7: **return** $\sqrt{2} \cdot 2^z$

---

$$\text{zeros}(k) = \text{largest } i \text{ s.t. } 2^i \mid k$$

```
0 b 1 1 1 0 1 1 0 0 0 1 0
```
$$\text{zeros}(k)$$

Space $s = O(\log n)$     (hash function)

$\quad\quad + O(\log\log n)$     (storing z)

$\quad\quad = O(\log n)$

1: Pick $h: [n] \to [n]$ from a strongly universal hashing family
2: $z \leftarrow 0$
3: **for all** $1 \le i \le m$ **do**
4:     Get item $a_i \in [n]$
5:     **if** zeros$(h(a_i)) \ge z$ **then**
6:       $z \leftarrow$ zeros$(h(a_i))$
7: **return** $\sqrt{2} \cdot 2^z$

Intuition:    "Since $h$ behaves like a random function"

→ $d$ values $h(j_1), \ldots, h(j_d)$ uniformly distributed in $[n]$

→ for each, proba to have at least $r$ trailing zeroes in binary is $\frac{1}{2} \cdots \frac{1}{2} = \frac{1}{2^r}$

→ for $r \approx \log_2 d$, by "union bound" we "should" have at least one hash with $r$ trailing zeroes (w. cst proba)

→ for $r \gg \log_2 d$, by union bound very unlikely to have any hash with $r$ trailing zeroes

Analysis

Let

$$Y_n = \sum_{j: b_j > 0} \mathbb{1}_{\text{zeros}(h(j)) \geq n} \quad , \text{ for } n \geq 0$$

```
1: Pick h: [n] → [n] from a strongly universal hashing family
2: z ← 0
3: for all 1 ≤ i ≤ m do
4:    Get item a_i ∈ [n]
5:    if zeros(h(a_i)) ≥ z then
6:       z ← zeros(h(a_i))
7: return √2 · 2^z
```

Observations:

① $\boxed{Y_n \geq 1 \quad \text{iff} \quad z \geq n}$  (for all $n \geq 0$)

② $\mathbb{E}[Y_n] = \sum_{j: b_j > 0} \Pr[\text{zeros}(h(j)) \geq n] \underset{\underset{h(j) \text{ is uniformly distributed}}{\uparrow}}{=} \sum_{j: b_j > 0} \frac{1}{2^n} \underset{\underset{d \text{ terms}}{\uparrow}}{=} \frac{d}{2^n}$

③ $\text{Var } Y_n \underset{\underset{\substack{\text{pairwise} \\ \text{independence}}}{\uparrow}}{=} \sum_{j: b_j > 0} \text{Var } \mathbb{1}_{\text{zeros}(h(j)) \geq n}$

$$\leq \sum_{j: b_j > 0} \frac{1}{2^n} = \frac{d}{2^n}$$

So

$$\mathbb{E}[Y_n] \leq \frac{d}{2^n} \quad, \quad \mathrm{Var}[Y_n] \leq \frac{2^n}{d} \quad \text{for all } n \geq 0$$

```
1: Pick h: [n] → [n] from a strongly universal hashing family
2: z ← 0
3: for all 1 ≤ i ≤ m do
4:     Get item a_i ∈ [n]
5:     if zeros(h(a_i)) ≥ z then
6:         z ← zeros(h(a_i))
7: return √2 · 2^z
```

④ Markov!

$$\Pr[z \geq n] = \Pr[Y_n \geq 1] \leq \mathbb{E}[Y_n] = \frac{d}{2^n} \quad ⊛$$

Chebyshev!

$$\Pr[z \leq n] = \Pr[Y_{n+1} = 0] \leq \frac{2^{n+1}}{d} \quad ⊛$$

⑤ Conclude:

$$C = 3\sqrt{2}$$

$$\Pr[\hat{d} \geq Cd] = \Pr\left[2^z \geq \frac{C}{\sqrt{2}}d\right] \overset{⊛}{\leq} \frac{\sqrt{2}}{Cd} \cdot d = \frac{1}{3}$$

$$\Pr[\hat{d} \leq d/c] = \Pr\left[2^z \leq \frac{d}{\sqrt{2}\,c}\right] \overset{⊛}{\leq} \frac{2d}{\sqrt{2}\,c} \cdot \frac{1}{d} = \frac{1}{3}$$

⑥ Amplify (?) Not a very good guarantee! Union bound naively gives $\Pr\left[\hat{d} \notin [d/c, Cd]\right] \leq \underbrace{2/3}_{\geq 1/2} \ldots$

How to amplify this?

"Carefully": median trick still applies, and works.

**Theorem 42.** *The (median trick version of the)* TIDEMARK *(AMS) algorithm is a* randomised *one-pass algorithm which, for any given parameter $\delta \in (0,1]$, provides an estimate $\hat{d}$ of the number $d$ of distinct elements of the stream such that, for some absolute constant $C > 0$,*

$$\Pr\left[\frac{1}{C} \cdot d \leq \hat{d} \leq C \cdot d\right] \geq 1 - \delta$$

*with space complexity*

$$s = O\left(\log n \cdot \log \frac{1}{\delta}\right).$$

Only issue: $C = \Theta(1)$.
We don't get $1 \pm \varepsilon$
for arbitrary $\varepsilon > 0$.

# Can we do better?

Yes.

# Third example: Distinct Elements, the BJKST algorithm (1/4)

**Input:** Parameter $\varepsilon \in (0, 1]$

1: Set $k \leftarrow O(\log^2 n / \varepsilon^4)$, $T \leftarrow \Theta(1/\varepsilon^2)$
2: Pick $h\colon [n] \rightarrow [n]$ from a strongly universal hashing family
3: Pick $g\colon [n] \rightarrow [k]$ from a strongly universal hashing family

4: $z \leftarrow 0$, $B \leftarrow \emptyset$
5: **for all** $1 \leq i \leq m$ **do**
6:      Get item $a_i \in [n]$
7:      **if** $\text{zeros}(h(a_i)) \geq z$ **then**
8:          $B \leftarrow B \cup \{(g(a_i), \text{zeros}(h(a_i)))\}$
9:          **while** $|B| \geq T$ **do**
10:             $z \leftarrow z + 1$
11:             Remove every $(a, b)$ with $b < z$ from $B$
12: **return** $|B| \cdot 2^z$

*Handwritten margin notes:*

Why this value of $k$?

Birthday paradox.
We want **no** collision in any of the $\log n$ buckets (with high proba).

$k = O((T \log n)^2)$

means each bucket has no collision under $g$ with proba
$\geq 1 - \frac{1}{10 \log n}$

+ union bound to get overall proba $1 - \frac{1}{10} = \frac{9}{10}$

# Third example: Distinct Elements, the BJKST algorithm (2/4)

Space:
- Hash functions: $h, g$ take space $O(\log n + \log k)$
$$= O(\log n + \log \frac{1}{\varepsilon})$$
- $z$: space $O(\log \log n)$
- $B$: space $T \cdot O(\log k + \log \log n) = O\left(\frac{1}{\varepsilon^2}\left(\log \log n + \log \frac{1}{\varepsilon}\right)\right)$

Total: $O\left(\log n + \dfrac{\log(1/\varepsilon) + \log \log n}{\varepsilon^2}\right)$

**Input:** Parameter $\varepsilon \in (0,1]$
1: Set $k \leftarrow O(\log^2 n/\varepsilon^4)$, $T \leftarrow \Theta(1/\varepsilon^2)$
2: Pick $h: [n] \to [n]$ from a strongly universal hashing family
3: Pick $g: [n] \to [k]$ from a strongly universal hashing family

4: $z \leftarrow 0$, $B \leftarrow \emptyset$
5: **for all** $1 \leq i \leq m$ **do**
6:     Get item $a_i \in [n]$
7:     **if** $zeros(h(a_i)) \geq z$ **then**
8:         $B \leftarrow B \cup \{(g(a_i), zeros(h(a_i)))\}$
9:         **while** $|B| \geq T$ **do**
10:             $z \leftarrow z+1$
11:             Remove every $(a, b)$ with $b < z$ from $B$

12: **return** $|B| \cdot 2^z$

Assumption: no collisions via hashing. Our setting of $k$ ensures this is true with high proba. ($\geq 9/10$), so we can just add $1/10$ of failure proba at the end (union bound) to account for it.

Analysis:

As before, $\mathbb{E}[Y_n] = \dfrac{d}{2^n}$, $Var[Y_n] \leq \dfrac{2^n}{d}$ for all $n \geq 0$.

We return
$$Y_z 2^z$$

on line 12, so we fail when $|Y_z 2^z - d| > \varepsilon d$.

# Third example: Distinct Elements, the BJKST algorithm (3/4)

$$\Pr[\text{fail}] = \Pr\left[\,|Y_Z 2^Z - d| > \varepsilon d\,\right] = \sum_{\eta=1}^{\log n} \Pr\left[\,|Y_\eta - \frac{d}{2^\eta}| > \varepsilon \frac{d}{2^\eta}\,,\ Z = \eta\,\right]$$

$$\leq \sum_{\eta=1}^{\log n} \min\left(\Pr\left[\,|Y_\eta - \frac{d}{2^\eta}| > \varepsilon \frac{d}{2^\eta}\,\right],\ \Pr[Z = \eta]\right)$$

$$\leq \sum_{\eta=1}^{s-1} \Pr\left[\,|Y_\eta - \frac{d}{2^\eta}| > \varepsilon \frac{d}{2^\eta}\,\right] + \sum_{\eta=s}^{\log n} \Pr[Z = \eta]$$  ✪

$s$ is a param. to choose

$$= \Pr[Z \geq s]$$
$$= \Pr[Y_{s-1} \geq T]$$

$$\leq \sum_{\eta=1}^{s-1} \frac{2^\eta}{\varepsilon^2 d} + \frac{d}{2^{s-1} T} \quad \text{Markov!}$$

Chebyshev

$$\boxed{\mathbb{E}[Y_\eta] = \frac{d}{2^\eta} \\ \mathrm{Var}[Y_\eta] \leq \frac{2^\eta}{d}}$$

setting of $T$

$$\leq \frac{2^s}{\varepsilon^2 d} + \frac{\varepsilon^2 d}{100 \cdot 2^s}$$

Balance the two terms:
$$2^s \approx \varepsilon^2 d$$

$$\leq \frac{1}{10}$$

and we are done!

↑ If we had the $s$-th bucket in play, the $(s-1)^{th}$ must have reached capacity before

## Input: Parameter $\varepsilon \in (0,1]$

1: Set $k \leftarrow O(\log^2 n / \varepsilon^4)$, $T \leftarrow \Theta(1/\varepsilon^2)$
2: Pick $h : [n] \rightarrow [n]$ from a strongly universal hashing family
3: Pick $g : [n] \rightarrow [k]$ from a strongly universal hashing family

4: $z \leftarrow 0$, $B \leftarrow \emptyset$
5: **for all** $1 \leq i \leq m$ **do**
6:     Get item $a_i \in [n]$
7:     **if** $\text{zeros}(h(a_i)) \geq z$ **then**
8:         $B \leftarrow B \cup \{(g(a_i), \text{zeros}(h(a_i)))\}$
9:         **while** $|B| \geq T$ **do**
10:            $z \leftarrow z + 1$
11:            Remove every $(a, b)$ with $b < z$ from $B$
12: **return** $|B| \cdot 2^z$

✪ Key idea: use one bound (+ Chebyshev) for some terms, the other (+ Markov) for the remaining ones.

**Theorem 43.** *The (median trick version of the)* BJKST *algorithm is a randomised one-pass algorithm which, for any given parameters $\varepsilon, \delta \in (0, 1]$, provides an estimate $\widehat{d}$ of the number $d$ of distinct elements of the stream such that, for some absolute constant $C > 0$,*

$$\Pr\left[ (1 - \varepsilon) \cdot d \le \widehat{d} \le (1 + \varepsilon)d \right] \ge 1 - \delta$$

*with space complexity*

$$s = O\left( \left( \log n + \frac{\log(1/\varepsilon) + \log \log n}{\varepsilon^2} \right) \cdot \log \frac{1}{\delta} \right).$$

# … Can we do better?

Yes.

(a little bit)

(But it's a much more complicated algorithm/analysis)