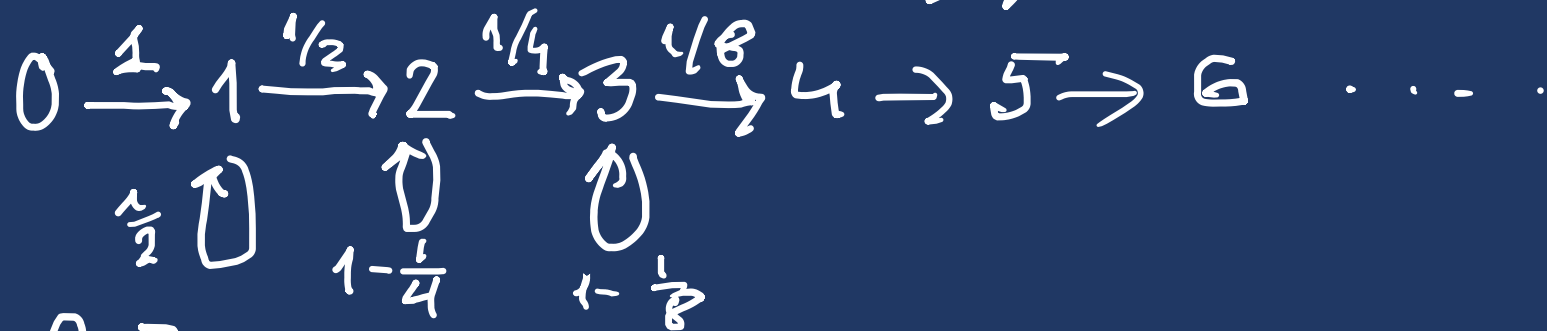


Morris Counter:

$$\begin{cases} C=0. \\ \text{onInc}(C) \{ \text{if}(\text{rand}() < (\frac{1}{2})^C) \{ C++ \} \} \end{cases}$$

$$\approx \underline{\underline{1.955}}$$



$$E[2^{C_n}] = 2$$

C_n = the value of C after n increments

$$E[2^{C_0}] = 2^0 = 1.$$

$$k < 0 : P[C_n = k] = 0$$

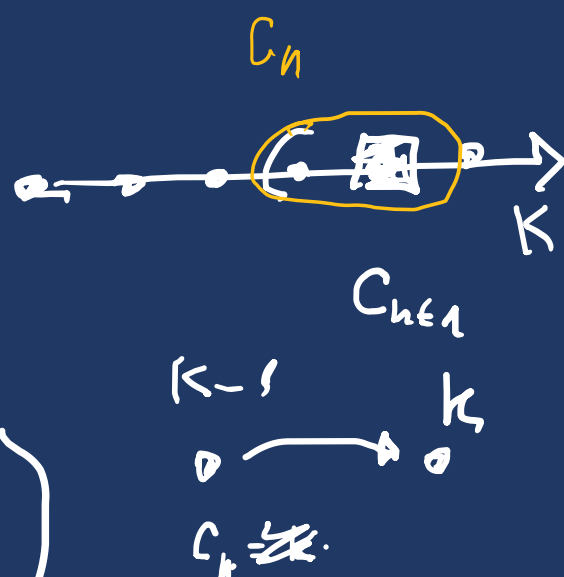
$$E[2^{C_{n+1}}] = \sum_k 2^k P[C_{n+1} = k] =$$

$$= \sum_k 2^k \left(P[C_{n+1} = k | C_n = k] P[C_n = k] + P[C_{n+1} = k | C_n = k-1] P[C_n = k-1] \right) =$$

$$= \sum_k 2^k \left(\left(1 - \frac{1}{2^k}\right) P[C_n = k] + \frac{1}{2^{k-1}} P[C_n = k-1] \right)$$

$$= \sum_k 2^k P[C_n = k] - \sum_k P[C_n = k] + \sum_k 2 \cdot P[C_n = k-1] =$$

$$= E[2^{C_n}] - 1 + 2 \sum_{l=0} P[C_n = l] = E[2^{C_n}] + 1.$$



$$\begin{cases} \cdot E[2^{C_0}] = 1 \\ \cdot E[2^{C_{n+1}}] = E[2^{C_n}] + 1 \end{cases}$$

$$\begin{aligned} E[2^{C_3}] &= E[2^{C_2}] + 1 = E[2^{C_1}] + 1 + 1 = \\ &= E[2^{C_0}] + \underbrace{1+1+1}_3 = 1 + 3 \end{aligned}$$

$$E[2^{C_n}] = n + 1$$

If we know C then

$$E[2^C - 1] = n$$

So we may take

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

$$\hat{n} = 2^C - 1$$

EXERCISE: count $\text{var}(2^{C_n})$,

$$\bullet \text{var}(2^{C_n}) = E[(2^{C_n})^2] - (E[2^{C_n}])^2$$

$$\text{var}(X) = E(X^2) - (E(X))^2$$

$$(2^{C_n})^2 = 2^{2 \cdot C_n} = 4^{C_n}$$

Try to count

$$E[4^{C_{n+1}}] = \sum_k 4^k P[C_{n+1} = k] = \dots$$

• Chebyshev inequality:

$$P[|2^{C_n} - E[2^{C_n}]| \geq \alpha \cdot E[2^{C_n}]] \leq \frac{\text{var}(2^{C_n})}{\alpha^2 (n+1)^2}$$

GEOMETRIC COUNTER

$\{a, a, b, a, b, b, c, a, \dots\} \longrightarrow \{a, b, c\}$

how many distinct elements
was in the stream (till now)

- use good hash function h

$\{a, a, b, a, b, b, \dots\} \rightsquigarrow \{h(a), h(a), h(b), \dots\}$

$a \neq b \longrightarrow P[h(a) = h(b)]$ is very small:

$h: \Sigma^* \longrightarrow \{0, 1\}^{128}$

$a \neq b \longrightarrow P[h(a) = h(b)] \approx \frac{1}{2^{128}}$ take h from
universal families

- $L = 0$;
- on Update(a) {

$$h(a) = (x_1 x_2 \dots x_l);$$

$$\rightarrow \omega(h(a)) = \min\{i : h(a)_i = 1\};$$

$$L = \max(L, \omega(h(a)));$$

}

~~a, b, c~~ a, b, c, c,

• $h(a) = (000\underline{1}000) \quad \therefore c = 4$

• $h(b) = (0\underline{1}10010) \quad \therefore c = 4$

— $h(a) = (0001\dots) \quad \therefore c = 4$

• $h(c) = (000\underline{0}10) \quad \therefore c = 5$

e.g. $l = 128$..

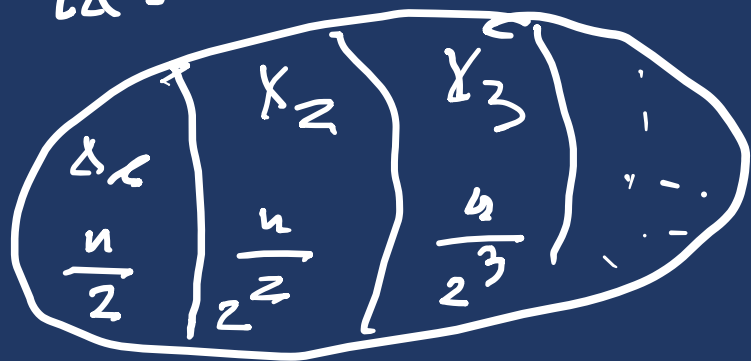
$$h(a) = (0, 0, 0, 1, 0, 0, \dots)$$

↑
4

~~$h(a) = 2$~~

first
one

$|X| = n$ X -set



$$\omega(a) = \text{mlu} \{k : h(a)_k = 1\}$$

$$X_k^0 = \{x \in X : \omega(x) = k\}$$

$$X_1 = \{x \in X : h(x) = "1*" \} \quad |X_1| \approx \frac{n}{2}$$

$$X_2 = \{x \in X : h(x) = "01*" \} \quad |X_2| \approx \frac{n}{2^2}$$

$$X_3 = \{x \in X : h(x) = "001*" \} \quad |X_3| \approx \frac{n}{2^3}$$

when $|X_k| \leq 1$? $\frac{n}{2^k} \leq 1$

$$n \leq 2^k$$

$$k \geq \lg_2 n !!!$$

L_n = the value of counter after observing n dist. elements.

$$P[L_n \leq k] = \left(1 - \frac{1}{2^k}\right)^n$$

- $X_{(1)}, \dots, X_{(n)}$ - indep. rand variables with $\text{Geo}\left(\frac{1}{2}\right)$

- $\tilde{L} = \max\{X_{(1)}, \dots, X_{(n)}\}$

$L_n \sim \tilde{L}$.

- $E[\tilde{L}] = \frac{1}{2} + \frac{1+n}{\ln 2} + \sigma_n$

quite difficult

$|\sigma_n| < 0.01$

fix n

$\approx \log_2 n$.



$\left(1 - \frac{1}{2^{\log_2 n}}\right)^n = \left(1 - \frac{1}{n}\right)^n \approx e^{-1}$

unintuitive.

IDEA : use $\hat{n} = 2^L$
as an estimate of n

$$E[L_n] \approx \log_2 n$$

Problem : precision

solution :

try to use many such
counters

$$L \approx \log_2 n$$
$$n \approx 2^L$$

L_1, L_2, \dots, L_m

each counter
has
hash function

try use

$$\frac{L_1 + L_2 + \dots + L_m}{m}$$

$$= \text{mean}(L_1, L_2, \dots, L_m)$$

many calculations

How to use only one hash function?

$$L = [L_0, L_1, \dots, L_{m-1}]$$

$\underbrace{\hspace{15em}}_m$

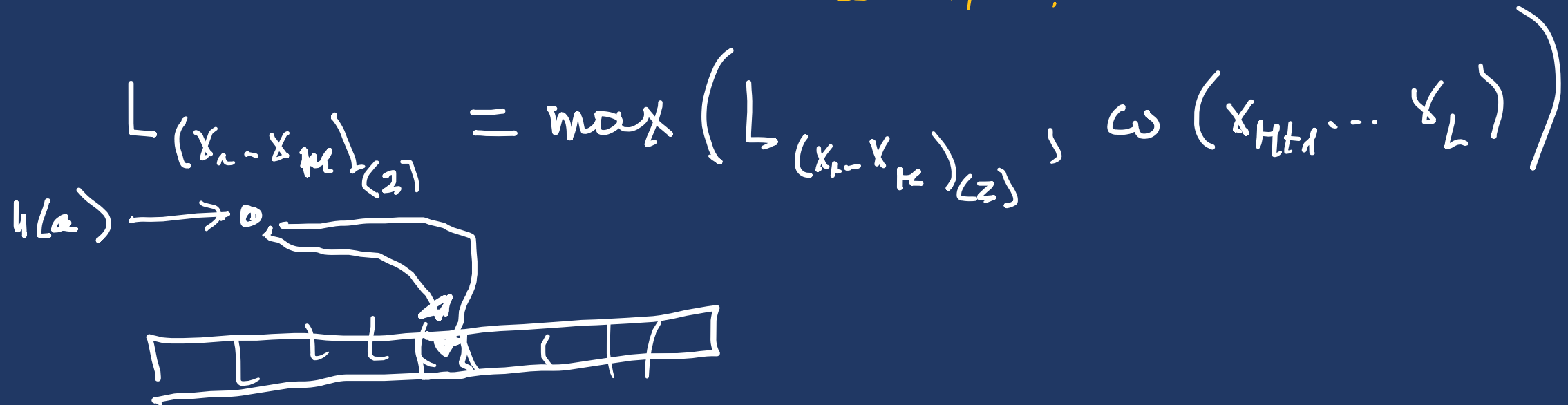
$$m \approx 64$$

$$\approx 128$$

$$\underline{\underline{m \approx 1024}}$$

$$x = h(a)$$

$$x = (x_1, x_2, \dots, x_L) = \underbrace{(x_1, \dots, x_M)}_{\text{coll. pos.}, m = 2^M} \parallel (x_{M+1}, \dots, x_L)$$



onUpdate(a) {

$$s = h(a)$$

$$s = s_1 \parallel s_2 ; \quad (|s_1| = M)$$

$$i = (s_1)_{(2)}$$

$$L[i] = \omega^*(s_2)$$

}

$$\omega^*(s) = \begin{cases} \min \{k : s_k = 1\} : (\exists i) (s_i = 1) \\ (L - M) + 1 : (\forall i) (s_i = 0) \end{cases}$$

~~W. 4013~~

$$h: \Sigma \rightarrow \{0, 1\}^L$$

HYPER-LOG-LOG

$$|X| = n$$



$$m$$

 \approx

$$\frac{n}{m}$$

$$\frac{n}{m} \approx 2^{L_i}$$

$$m \approx m \cdot 2^{L_i} \quad i = 0, 1, \dots, m-1$$

$$i = 0, \dots, m-1$$

$$\hat{m}_i \approx 2^{L_i}$$

$$\log_2(\log_2 n)$$

Try to use

$$\text{mean} (m 2^{L_0}, m 2^{L_1}, \dots, m 2^{L_{m-1}})$$

as an estimator of n

Means

$$x_1 \dots x_m > 0$$

$$\frac{1}{\frac{1}{x_1} + \dots + \frac{1}{x_m}} \leq \sqrt[m]{x_1 \dots x_m} \leq \frac{x_1 + \dots + x_m}{m}$$

harmonic mean geometric mean arithmetic mean

$$\frac{m}{\frac{1}{m \cdot 2^{L_0}} + \frac{1}{m \cdot 2^{L_1}} + \dots + \frac{1}{m \cdot 2^{L_{m-1}}}} = \frac{m^2}{\binom{1}{2}^{L_0} + \dots + \binom{1}{2}^{L_{m-1}}}$$

this is our estimate

HLL

on GetEstimate () {
 $Z = m^2 / \sum_{c=0}^{m-1} \left(\frac{1}{2} \right)^{L[c]}$;

return ($\alpha_m \cdot Z$)
 }

$$\alpha_m = \left(m \int_0^{\infty} \left(\log_2 \left(\frac{u+2}{u+1} \right) \right)^m du \right)^{-1}$$

HyperLogLog
Algorithm

m	α_m
16	0.673
32	0.69
64	0.709
≥ 128	$\frac{0.7213}{1 + \frac{1.074}{m}}$

FINAL IMPROVEMENTS



• count $z = |\{i: L_i = 0\}|$
IF $z \neq 0$

use linear counting, e.g.

return: $m \ln \frac{m}{z}$,

• if some $i: L_i >$

correction the $L_i \leftarrow -2^{64} \ln \left(1 - \frac{L_i}{2^{64}}\right)$

return $(\alpha_m \cdot z)$

EXER.
Implement
it

HYPER LOG-LOG

precision :

$$\frac{1.04}{\sqrt{m}}$$

$$m = 124$$

:

$$\sigma \approx \frac{1}{10} = 10\%$$

$$m = 1624$$

:

$$\sigma \approx \underline{\underline{3\%}}$$