

STREAMING : generate random sample



unif. distributed sample

$$A \subseteq \{1 \dots N\} \quad |A| = k \quad P[S = A] = \frac{1}{\binom{N}{k}}$$

Basic alg: $k = 1$.

INIT: $p = 0$; $data = ucl$; $n = 0$;

on Get(x) { if (random() < $\frac{1}{n}$) { $p = n$; $data = x$ }

$n++$

Modifications

on Get(x) {

 n++;

 if (random() < $\frac{2}{n}$) {

 p = n;

 data = x

}

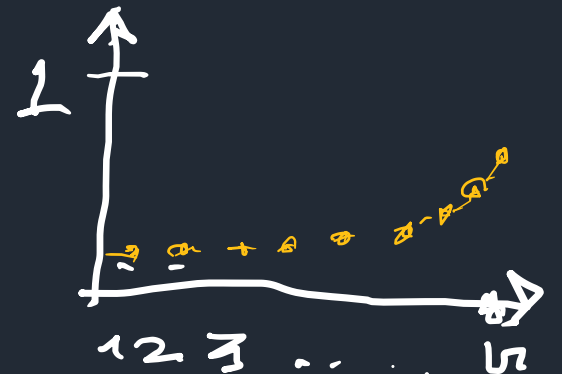
concentration $\frac{1}{n}$



(random() < $\frac{1}{n}$)

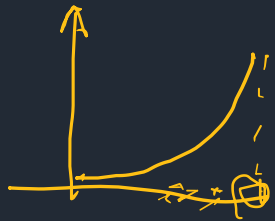


(random() < $\frac{2}{n}$)

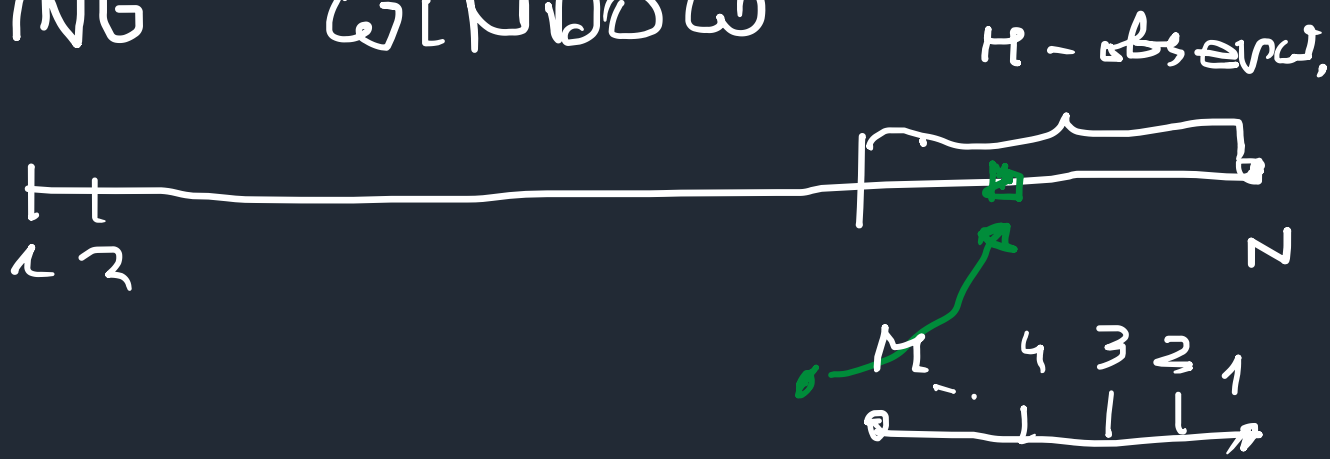


$< \frac{1}{n}$

$< \frac{1}{\sqrt{n}}$



SLIDING WINDOW



$(M = 10^9)$
 $N \gg 10^9$

p - rand. variable in $\{1, \dots, N\}$

• 2005: P. Indyk,
 Motwani

1) $1 \leq p \leq M$

2) on Get(x) {

• 2012: Burawerwan, ...

if ($p = M$) or (rand.dec.(p)) {

$p = 1;$

} data = x

COUNTING :

- we observe $(a_L)_{L=1..N}$
- calculate (estimate) :

$$\left\{ \{ a_L : L=1..N \} \right\}$$

(EX) $aabbaabbcccbbaa$ = (*)

$$\left\{ \{ a, b, c, d \} \right\} = 4$$

naive solution : (*) \Rightarrow $a a a a a b b b b c c d d$

(in memory)

sort 1 1 1 1 1 2 2 2 2 3 3 4 4

time: $O(n \log n)$

$\Rightarrow 4$

DO NOT LOOK AT a 's

CONSIDER $h(a)$ h - good hash function

$\{a_1, a_2, a_3, \dots, a_n\}$

$\{a_1, a_2, a_3, \dots, a_n\}$

$\{h(a_1), h(a_2), \dots, h(a_n)\}$

$i < j : a_i = a_j \rightarrow h(a_i) = h(a_j)$

suppose: $a \neq b \rightarrow h(a) \neq h(b)$ with high prob.
collision

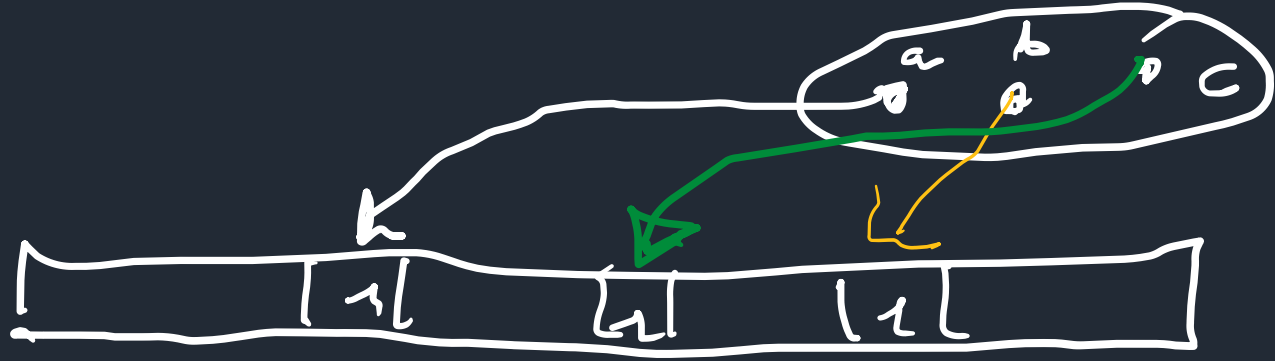
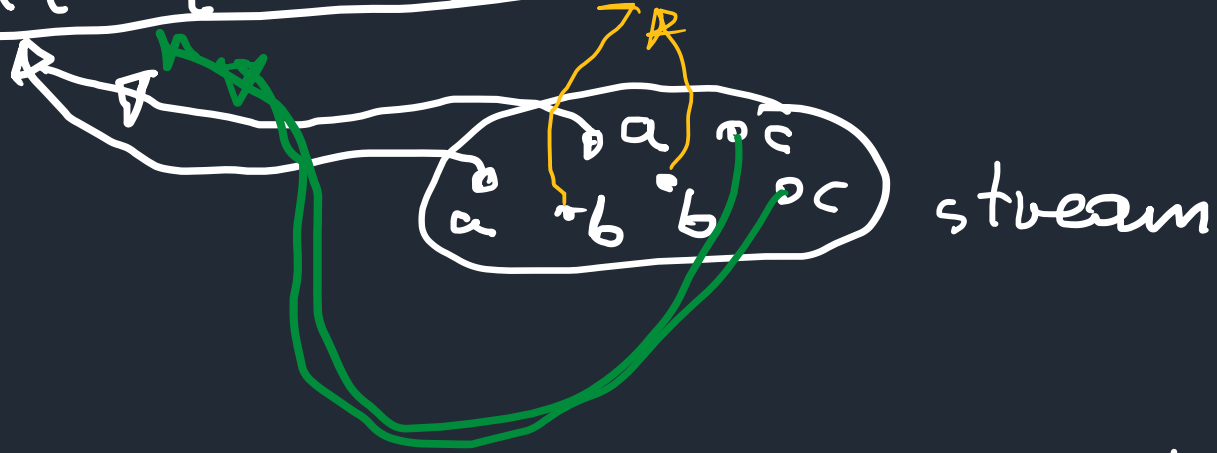
LINEAR COUNTING.

We have a table $X[1 \dots N]$
of bits (in memory).

~~Initially~~ $X = (0, 0, \dots, 0)$
we have a hash function $h: \Omega^* \rightarrow \{1, \dots, N\}$

on Get(a) {
 $i = h(a)$;
 $X[i] = 1$;
}

\equiv { on Get(a) {
 $X[h(a)] = 1$;
}



$$\left\{ \begin{array}{l} h(a), h(b) \\ h(c) \end{array} \right\} - \text{indep}$$



$$P[X[2] = 0] = \left(\frac{N-1}{N}\right)^m = \left(1 - \frac{1}{N}\right)^m$$

$$L = \sum_{i=1}^N \mathbb{I}[X[i] = 0] \quad (= \text{number of } i \text{ s.t. } X[i] = 0)$$

$$\underline{\underline{E[L]}} = \sum_{i=1}^N E(\mathbb{I}[X[i] = 0]) = \sum_{i=1}^N \Pr[X[i] = 0] = \underline{\underline{N\left(1 - \frac{1}{N}\right)^m}}$$

$$E[L] = N \left(1 - \frac{1}{N}\right)^m$$

m - number
of distinct
elements

Method of first moment

u = number of empty cells.

idea: $u \approx E[L]$

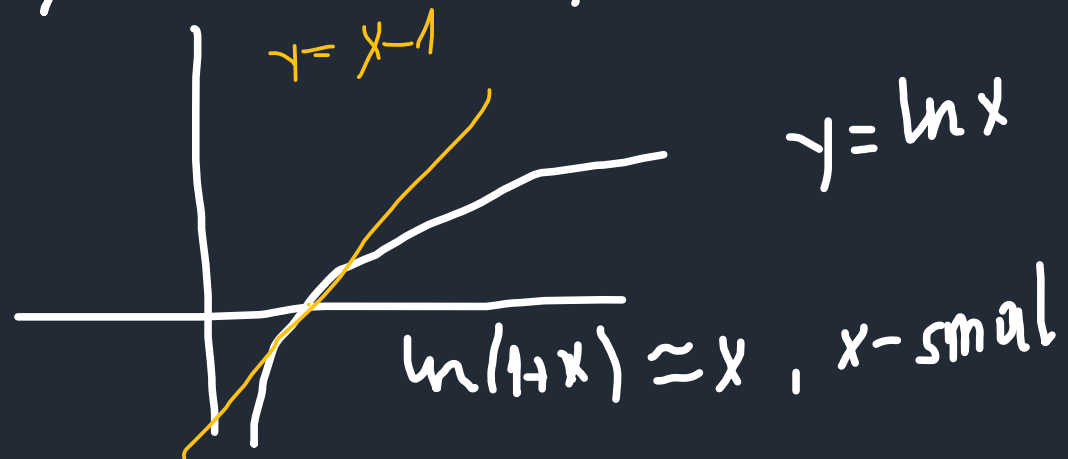
1 0 1 0 1 1 0 0 0 1

$$u = N \left(1 - \frac{1}{N}\right)^m$$

$$\frac{u}{N} = \left(1 - \frac{1}{N}\right)^m$$

$$\ln\left(\frac{u}{N}\right) = m \ln\left(1 - \frac{1}{N}\right)$$

$$m \approx \frac{\ln\left(\frac{u}{N}\right)}{\ln\left(1 - \frac{1}{N}\right)}$$



$$\ln\left(1 - \frac{u}{N}\right) \approx -\frac{u}{N}$$

$$\hat{m} = \frac{\ln\left(\frac{u}{N}\right)}{-\frac{1}{N}} = -N \ln\left(\frac{u}{N}\right) = N \ln\left(\frac{N}{u}\right)$$

of estimator
 \hat{m}

$$\hat{m} = N \ln\left(\frac{N}{u}\right)$$

!!!

u_1, u_2, u_3, \dots



m balls

Problem: $u = 0 \Rightarrow \hat{m} = \infty$

Coupon collector problem

$$m \approx n \ln n$$

$$m \ll n \ln n$$

$m < N \rightarrow$ no problem

if $m = N$: $E[L] = N \left(1 - \frac{1}{N}\right)^N \approx$
 $\approx N \cdot e^{-1} \approx \frac{2}{3}$

Example: $N = 100$
 $\ln N \approx 4.6$
 $N \cdot \ln N \approx 460$

$m \leq 250$; $n \leq 300$ \Rightarrow

precision \hat{m} $\begin{matrix} \uparrow \\ 0 \end{matrix}$

EXERCISES:

① implement LC,

② estimate

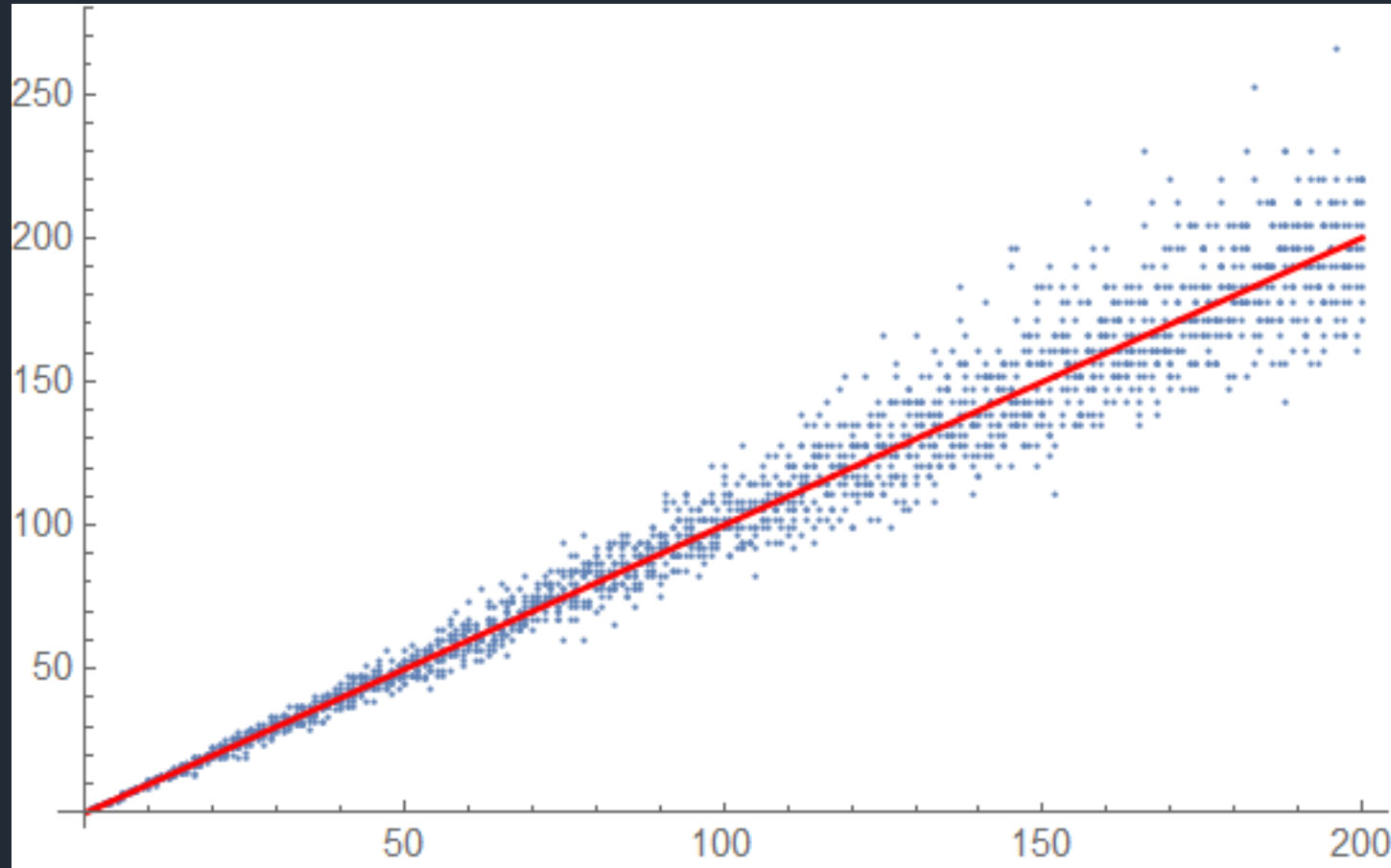
$$\text{var}(L) = E[L^2] - E[L]^2$$

$$\frac{\text{var}(L)}{E[L]^2}$$

thebyssen
weq.

③

Precision of LC for $N=100$ and $m=1,2,\dots,200$



(10 simulations for each k)

LC : good

$$m \leq 10^{10}$$

$$10^{10} = 10 \cdot \underbrace{10^9}_{1.5B}$$

we need

$\approx 1.5B$ Memory

for X .

{ collisions
in Genevix :

$$10^{15} \approx m$$

to buy for LC.

PROBABILISTIC COUNTERS

• classical counter

• $n = 0$;

• on Get(x) { u+t }

$n \leq N$

$$n = (a_k a_{k-1} \dots a_2 a_1 a_0)_2 = a_i \in \{0,1\}$$

$$= a_0 + a_1 \cdot 2 + a_2 \cdot 2^2 + \dots + a_k \cdot 2^k$$

$$\leq 1 + 2 + \dots + 2^k \leq 2^{k+1} - 1 \leq N$$

$$2^{k+1} \leq N + 1$$

$$k \leq \log_2(N+1) - 1$$

we need $\lceil \log_2(N) \rceil$ bits

Supposes that we do not need total precision!

$$n = 23645 \approx 23600 \approx 23700$$

$$n \approx 23500$$

MORRIS COUNTER

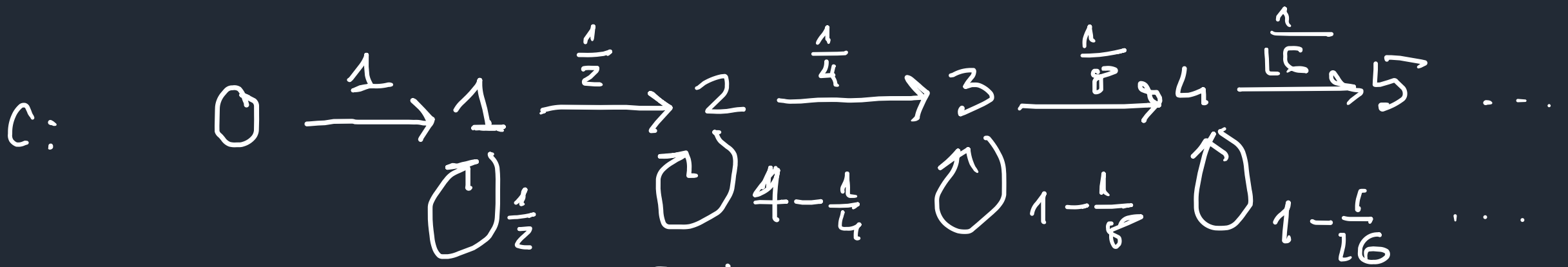
$C = 0$;

on Get (1C) {

if (random() < $\left(\frac{1}{2}\right)^C$) {

$C++$

}



if $(\text{round } C) < \left(\frac{1}{2}\right)^C \} C++$

Morris counter

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

C_n = the value of C after n increments.

we do it n times

$$2^{C_n} \approx n$$

$$C_n \approx \log_2 n$$

How many bits we need

$$\approx \log_2 (\log_2 n) \quad \square$$