

STREAMING

Sampling



```
on Get(x) {  
  n++;  
  if (random() < 1/n) {  
    p = n;  
    data = x  
  }  
}
```

otrzymujemy zmienne
losowe o
rozkr. jedu.

(sample with replacemnts)



MOZLIWE ZMIANY:



```
if (random() < 2/n) { ... }  
if (random() < 1/sqrt(n)) { ... }
```

SLIDING WINDOW



N = liczba ostatnich obserwacji

$$\text{rng}(X) \in \{1, \dots, N\}$$

rozkł. jednost.

± • 2003 : P. Indyk, Motwanił, ...

± • 2012 : Brauerman ← prosty

± • 2021 : D. Bójko, M. K., J. C.

COUNTING

string s : $(a_i)_{i=1..N}$

ALPHABET: $|\{a_i : i \in 1..N\}| = 2$

(P) $aabcaaaab$; $|\{a, b, c\}| = 3$

ΚΑΛΩΝΕ; sort: $aaaaabbb$

id From	id To
- -	- -
- -	- -

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

$|\{h(a_i) : i=1..n\}| \leq |\{a_i : i=1..n\}|$

LINEAR COUNTING

Many hash function $h: \Sigma^k \rightarrow \{1, \dots, N\}$

Many table bits $X[1 \dots N]$

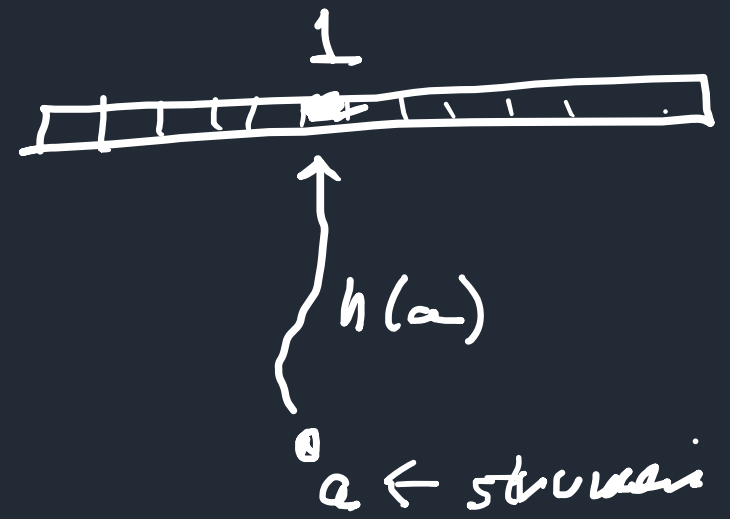
INIT: $X \leftarrow (0, 0, \dots, 0)$.

on $get(a)$ {

$X[h(a)] = 1$

}

precyfaligun m elementow.



$u, 42$



$a \neq b : h(a) \text{ vice } h(b)$

$$P[u_i = a] =$$



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

$$L = \sum_{l=1}^m \mathbb{I}[u_l = 0]$$

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

$u \leftarrow$ liczba pustych cerk

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



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

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

$$\begin{aligned} &\approx -N \ln\left(\frac{u}{N}\right) \\ &= N \ln \frac{N}{u} \end{aligned}$$

$$\hat{m} = N \ln \frac{N}{u}$$

Esty matrica liaboj m ,

$$\begin{aligned} \ln \frac{1}{1-x} &= x + \frac{1}{2}x^2 + \frac{1}{3}x^3 \\ &\Rightarrow \sum_{k \geq 1} \frac{1}{k} x^k \end{aligned}$$

$$\begin{aligned} \ln \frac{1}{1 - \frac{1}{N}} &\approx \frac{1}{N} \\ &= \end{aligned}$$

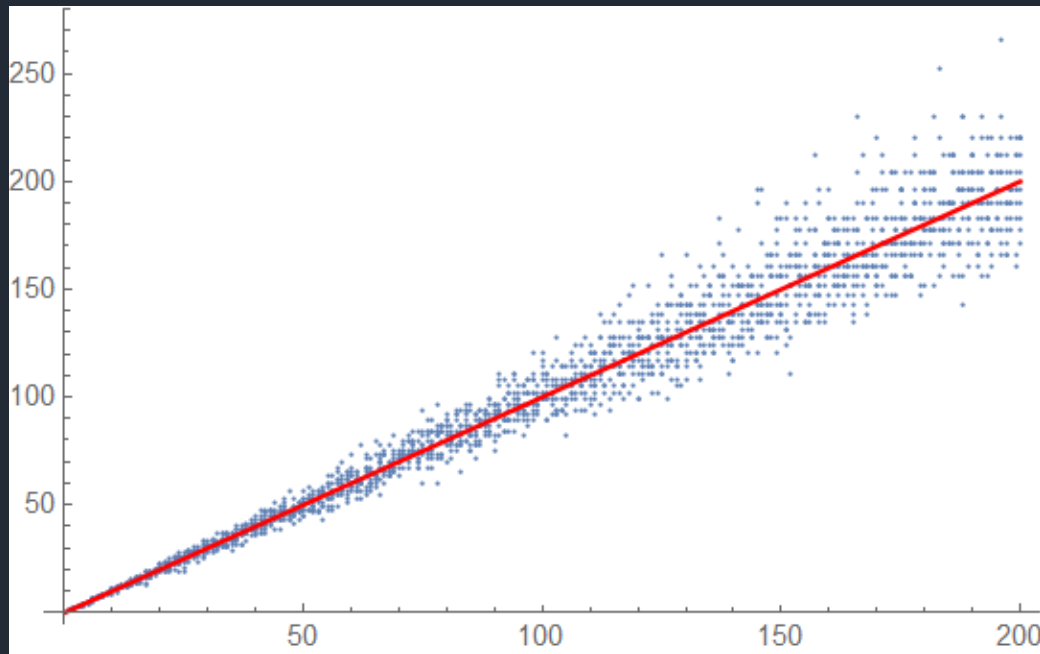
$$- \ln\left(1 - \frac{1}{N}\right)$$

PROBLEM: $u = 0$
 $? m < N \rightarrow NIE$
 Q&P: $u = 0 \rightarrow m = N \ln N$

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

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

$$P[|L - E[L]| \geq \alpha E[L]} \leq \frac{\text{var}(L)}{\alpha^2 (E[L])^2}.$$



simulations

$N = 100$;

$k = 1, 2, \dots, 200$

(10. simulations for each k)

Ograniczenie LIN. COUNTING : $m \ll N \ln N$

LICZNIK PROBABILISTYCZNY

• Licznik ; $\begin{cases} \text{init} : n = 0 ; \\ \text{onGet}(x) \{ n + z^x \} \end{cases}$

Ilo bitów potrzebuję do zapisania $n \in \{1, \dots, N\}$

$$L \approx \lg_2 N . \quad a_0 + a_1 2 + \dots + a_k 2^k \leq$$

• ≈ 1955 ; MORRIS IMB $\leq 1 + \dots + 2^k = 2^{k+1} - 1 \leq N$

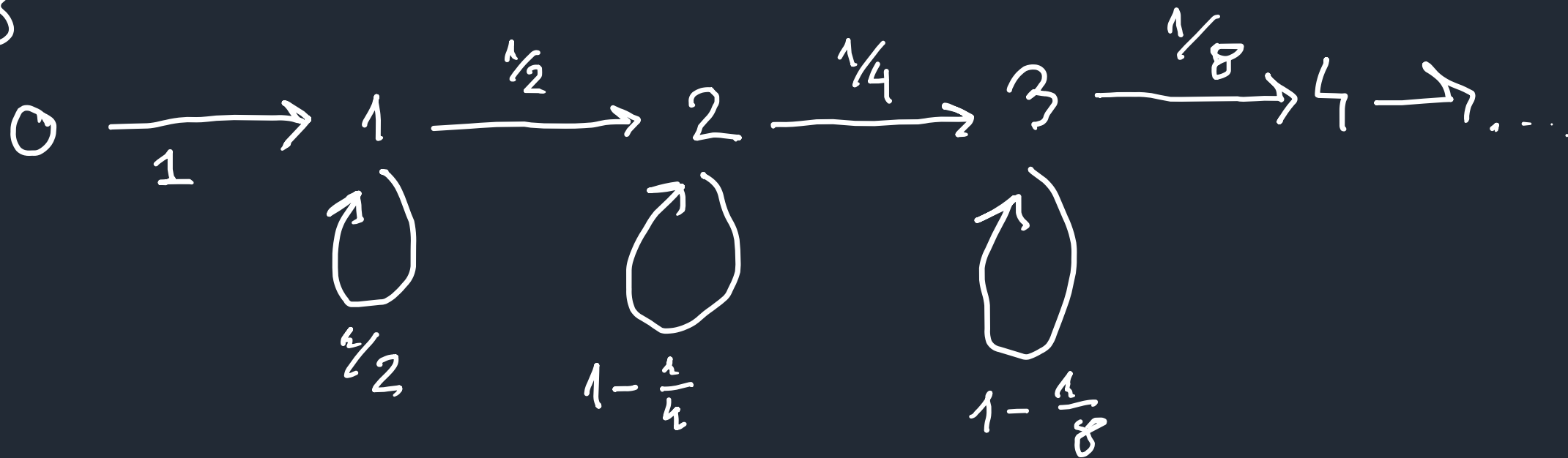
IDEA : nie muszą
liczyć dokładnie !

Liczuk Morrisa:

init: $C = 0$

onGet(x) {

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



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

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

$$E[2^{C_1}] = E[2^1] = 2,$$

$n \geq 1$

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

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

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

$$= \underbrace{\sum_k 2^k P[C_n = k]}_{E[2^{C_n}]} - \underbrace{\sum_k P[C_n = k]}_1 + 2 \underbrace{\sum_k P[C_n = k-1]}_1$$

C_n = wartość licznika po przecytniku n -elem.

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

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

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

C - wartość licznika :

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

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

$$2^c = n + 1$$

$$c \approx \log_2(n+1)$$

ile potrzebujemy bitów $n \rightarrow c$?

$$\log_2(\log_2(n+1))$$

