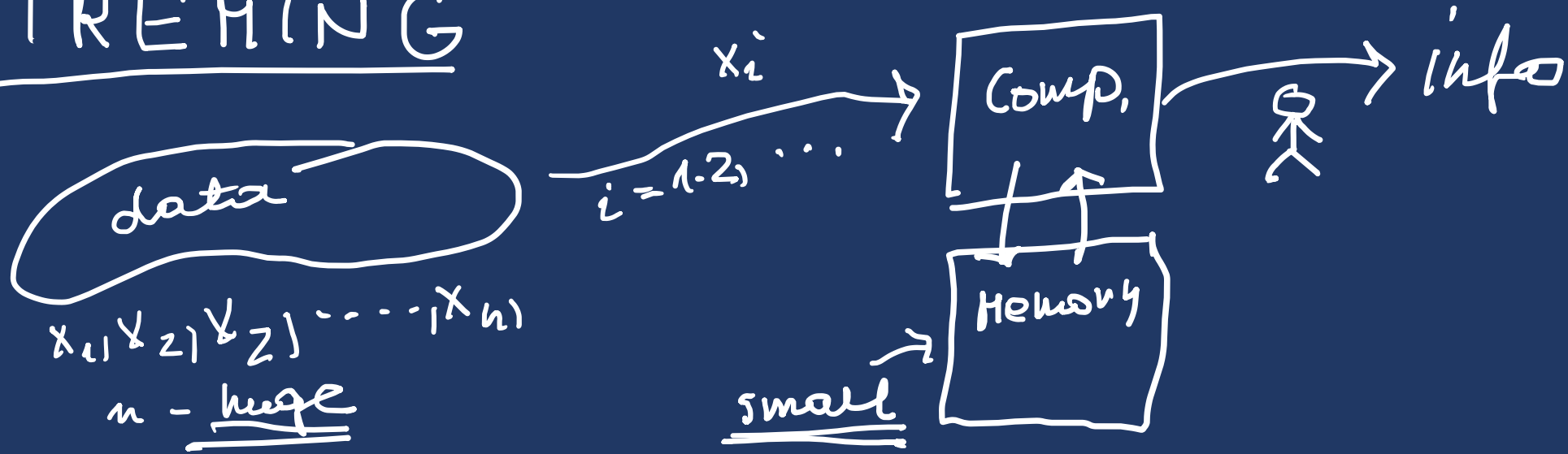
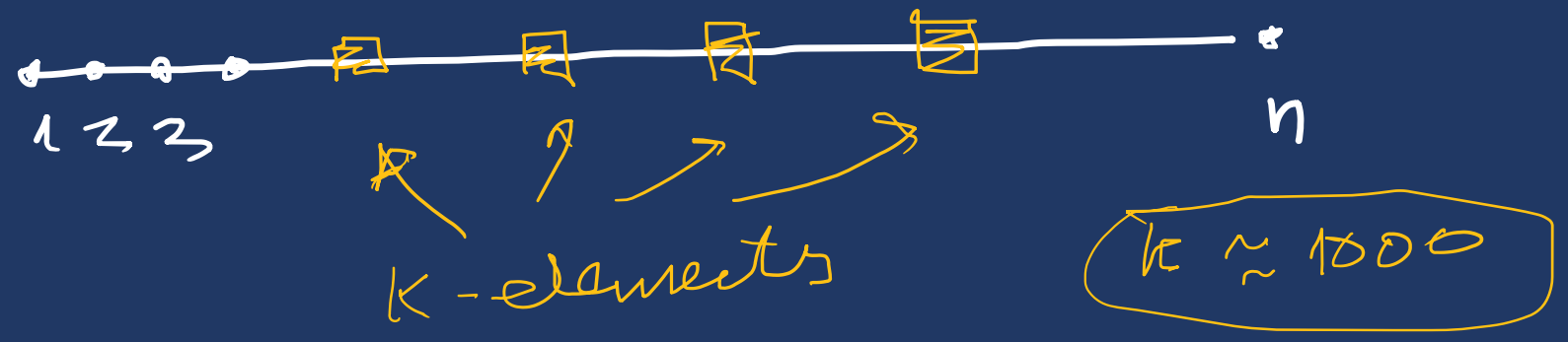


# STREAMING

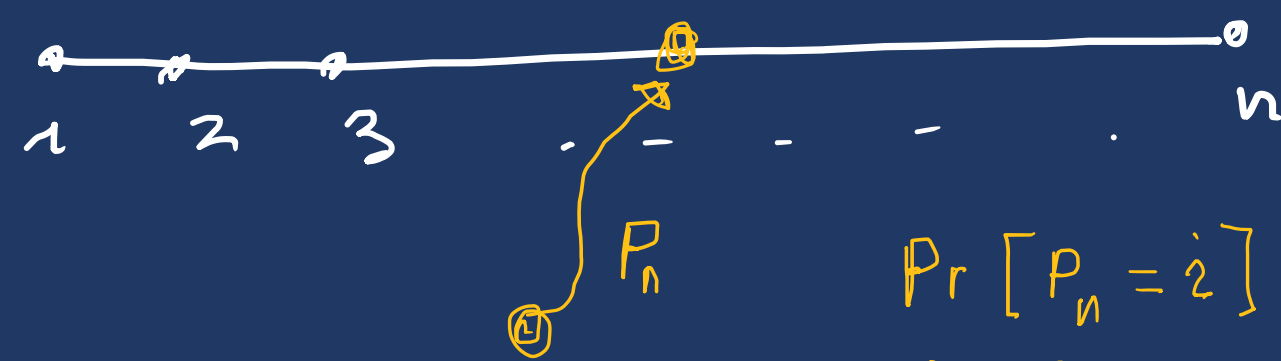


- Q : 1)  $n$  ?  
2)  $(\sum_{i=1}^n x_i / n)$   
3)  $(\sum_{i=1}^n x_i^2 / n)$  } trivial

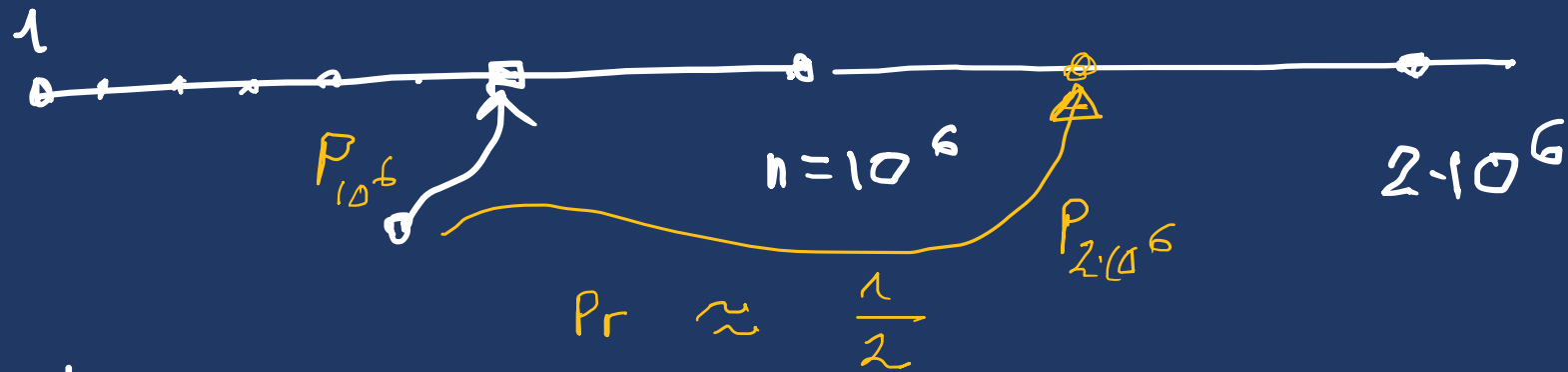
# SAMPLING :



sample : uniformly distributed



$$\Pr [P_n = i] = \frac{1}{n}$$
$$i \in \{1 \dots n\}$$



## Nitter's R-algorithm.

```

on Init ( ) {
  n = 0 ; // number of elem.
  p = 0 ; // pointer
  data = nil ; // obs. data by pointer
}

```

```

on Get ( ) {
  return ( p, data );
}

```

```

on Read ( x ) {
  n ++ ;
  if ( random (0,1) < 1/n ) {
    p = n ;
    data = x ;
  }
}

```

$P_n$  = the value of pointer after reading  $n$ -th element.

- $P_n \in \{1, \dots, n\}$

- GOAL ;  $(\forall n) (\forall i \in \{1, \dots, n\}) \left( \Pr [P_n = i] = \frac{1}{n} \right)$  (\*)

1)  $P_1 = 1$  : OK.

2) assume that (\*) holds at  $n$ ,  $C =$  "there was a change at  $n$ "

$$P[P_{n+1} = i] = P[P_{n+1} = i | C] \cdot \Pr[C] + P[P_{n+1} = i | \neg C] \cdot P[\neg C]$$

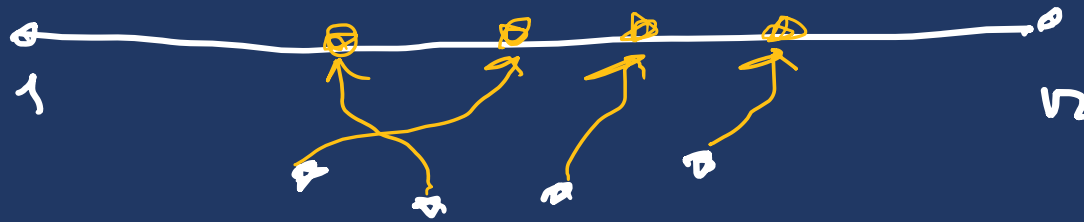
- $i \in \{1, \dots, n\}$  :  $P[P_{n+1} = i] = 0 \cdot \frac{1}{n+1} + P[P_n = i] \cdot \left(1 - \frac{1}{n+1}\right) =$   
 $= \frac{1}{n} \cdot \frac{n+1-1}{n+1} = \frac{1}{n+1}$

- $i = n+1$  :  $P[P_{n+1} = i] = 1 \cdot \frac{1}{n+1} + 0 \cdot \left(1 - \frac{1}{n+1}\right) = \frac{1}{n+1}$



- Build a sample of size  $k$   
use independently  $X_1, \dots, X_k$  from  $R$ -alg.

$$P[A_1 \cup \dots \cup A_k] \leq P(A_1) + \dots + P(A_k)$$



$$P\left[\bigwedge_{1 \leq i < j \leq k} (X_i \neq X_j)\right] = 1 - P\left[\bigvee_{1 \leq i < j \leq k} (X_i = X_j)\right] \geq$$

$$1 - \sum_{1 \leq i < j \leq k} P[X_i = X_j] = 1 - \binom{k}{2} \frac{1}{n} = 1 - \frac{k(k-1)}{2} \cdot \frac{1}{n} \approx$$

$$\approx 1 - \frac{k^2}{2n}$$

- if  $n \gg \frac{k^2}{2}$  then w.h.p.  $X_i \neq X_j$  for all  $1 \leq i < j \leq k$

Birthday Paradox  $\approx \sqrt{n}$ .

• Different version

we keep in memory array  $[x_1, \dots, x_k] = X$

$$x_i = (p_i, \text{data}_i)$$

$$(\forall n) (\forall A \subseteq \{1..n\}) (|A|=k \rightarrow P(\{p_1, \dots, p_k\} = A) = \frac{1}{\binom{n}{k}})$$

• read  $x_1, \dots, x_k$ ; put

$$[(1, x_1), \dots, (k, x_k)]$$

• if  $n > k$ :

$$\text{if } (\text{random}(\text{bool}) < \frac{k}{n}) \{$$

$$i = \text{random} \{1, \dots, k\}$$

$$X[i] = (n, x)$$

}

THIS IS CORRECT

EXERCISE

$$\binom{n}{k} \quad \binom{n+1}{k}$$



modification with prob  $\frac{k}{n}$ .

change  $\rightarrow$   $M$   $M+1$   $M+2$   $M+3$   $M+4$



$\frac{1}{M}$

$L \leftarrow$  s.t. next change is at  $M+L$

$$P[L \geq 1] = 1$$

$$P[L \geq 2] = 1 - \frac{1}{M+1} = \frac{M}{M+1}$$

$$P[L \geq 2] = P[L \geq 3] = \left(1 - \frac{1}{M+1}\right) \cdot \left(1 - \frac{1}{M+2}\right) = \frac{M}{M+1} \cdot \frac{M+1}{M+2}$$

$$P[L \geq k] = \frac{M}{\cancel{M+1}} \cdot \frac{\cancel{M+1}}{\cancel{M+2}} \cdot \dots \cdot \frac{M+(k-1)}{M+k} = \frac{M}{M+k}$$

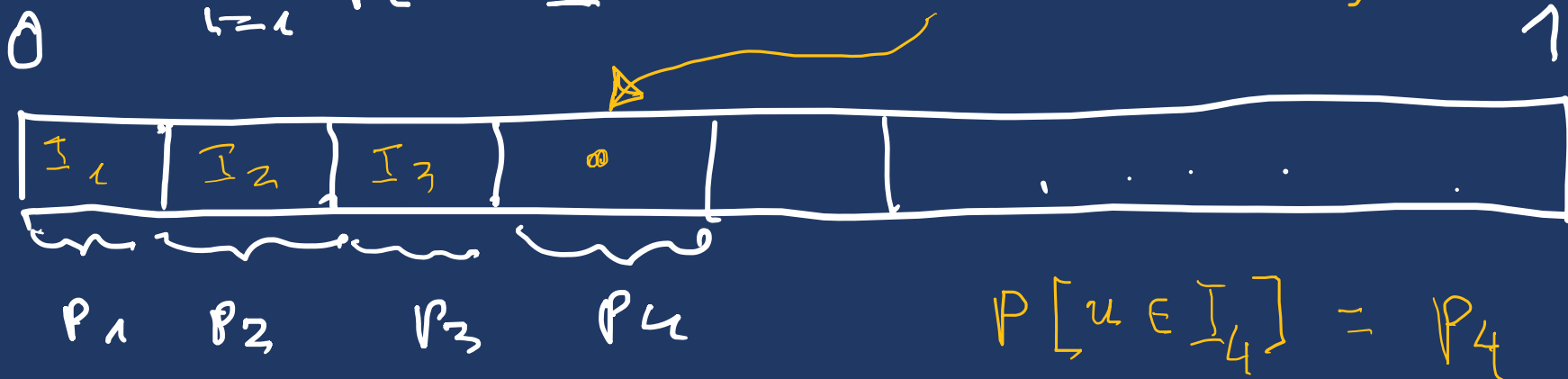
$$P[L \leq k] = 1 - \frac{M}{M+k} = \frac{k}{M+k} \xrightarrow{k \rightarrow \infty} 1$$

$$P[L \leq k] = \frac{k}{M+k}$$

$$P_{i0} = P[L=i] = \frac{i}{M+k} - \frac{(i-1)}{M+k-1}, \quad i=1,2,3,\dots$$

$$\sum_{i=1}^{\infty} P_i = 1$$

$$u = \text{rand}(0,1)$$



$$P[L \leq k-1] = p_1 + p_2 + \dots + p_{k-1} < u \leq p_1 + \dots + p_k = P[L \leq k]$$

$u = \text{rand}(0,1)$ ; find  $k$  s.t.  $(*)$ ; return  $k$ .

$$\text{find min } k : u \leq \frac{k}{M+k}$$



$$u \leq \frac{k}{M+k}$$

$$Mu + k \cdot u \leq k$$

$$k - k \cdot u \geq M \cdot u$$

minimal  $\rightarrow k \geq \frac{M \cdot u}{1 - u}$

$$k = \left\lceil \frac{M \cdot u}{1 - u} \right\rceil$$

$$\left\{ \begin{array}{l} u = \text{rand}(0, 1) \\ \text{return } \text{ceiling}((M \cdot u) / \\ (1 - u)); \end{array} \right.$$

init:  $n = 0$ ;  $p = 0$ ;  $data = \text{nil}$ ;  $c = \text{nil}$ ;  $L = 1$ ;

on Read ( $x$ ) {

$n++$ ;  $c++$ ;

if ( $c == L$ ) {

$p = n$ ;

$data = x$ ;

$u = \text{random}(0, 1)$ ;

$L = \text{ceiling}((n \cdot u) / (1 - u))$ ;

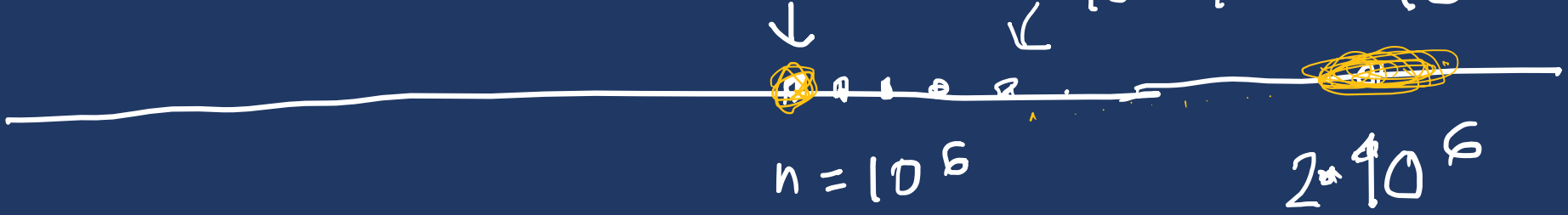
$c = 0$ ;

}

good version  
of R-algorithm

change

$$\frac{1}{10^6 + 4} \approx \frac{1}{10^6}$$



$$L \sim \text{Geo}\left(\frac{1}{10^6}\right)$$

next change  $\approx 2 \cdot 10^6$

