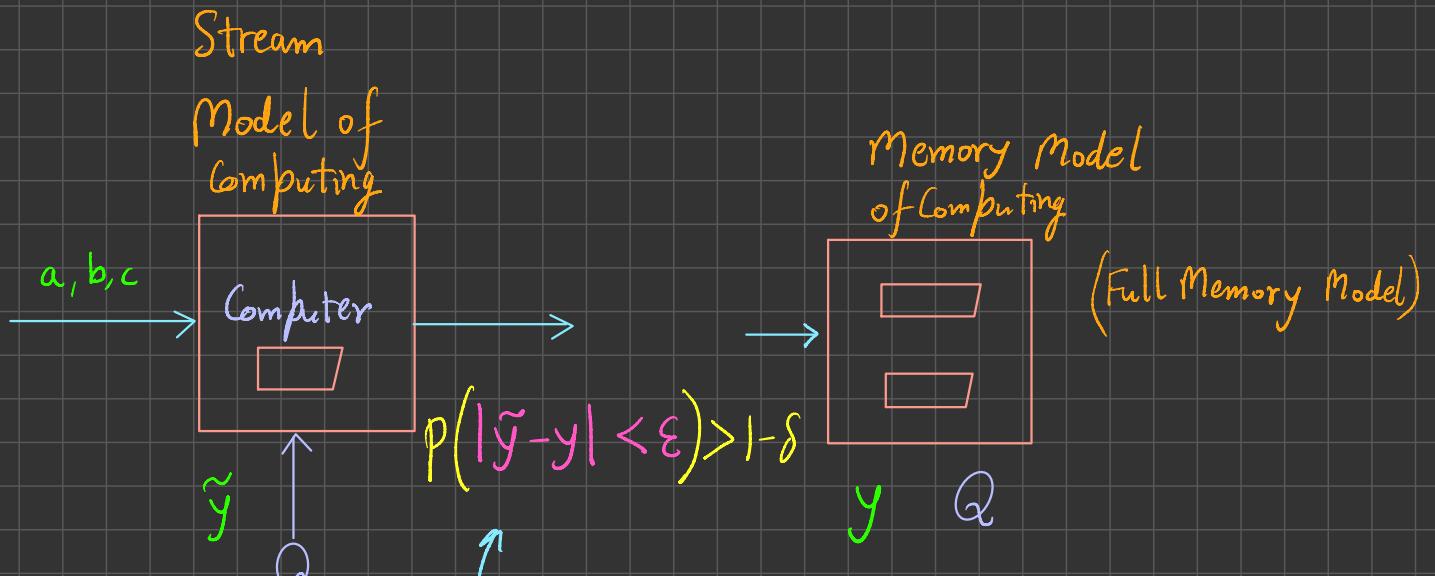


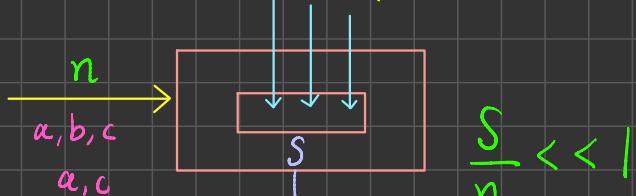
# Streaming data algorithms:



So look for lowering the prediction error.

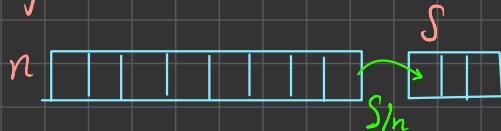
## Sampling

→ each element is equiprobable

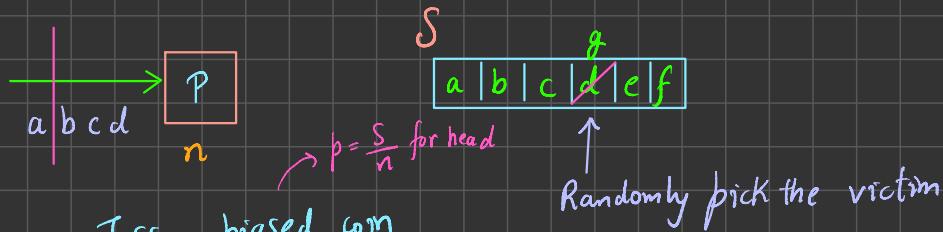


limited memory,  
so sample is a representation of all n elements

Faithful - Unbiased

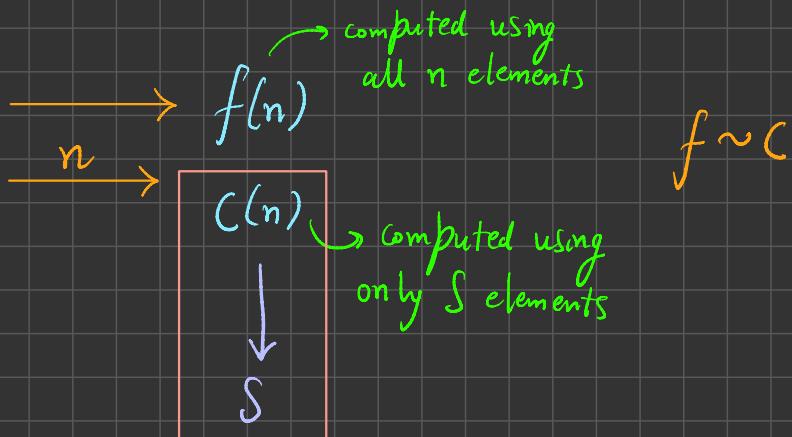


# Reservoir Sampling

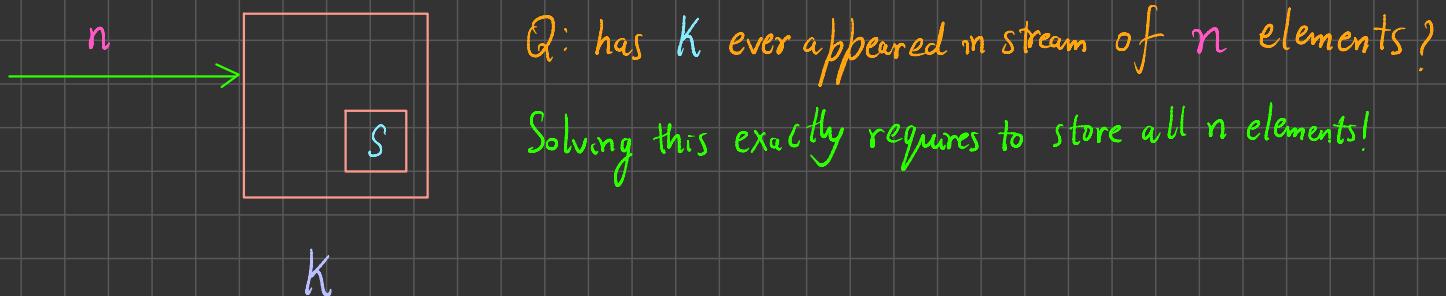


Toss a biased coin if head, then put the new element in  $S$ , else replace the element with  $(S+1)^{st}$  if heads on coin toss.

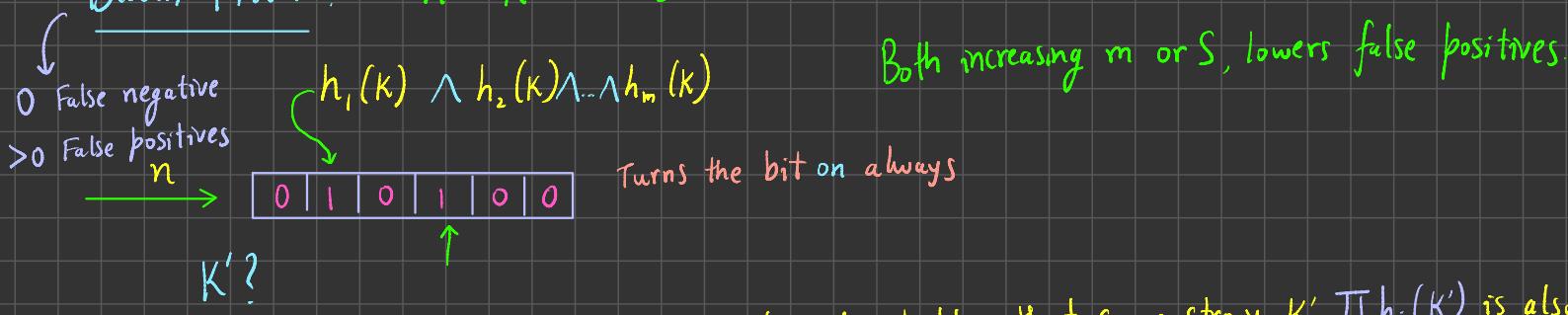
General idea: Sketches



Membership:



Bloom Filter:  $h: K \rightarrow i \in S$



If  $K'$  ever appeared, then detected, but it might happen that some stray  $K' \cap h_i(K')$  is also 1, which results in false positives

## Count distinct

$a \in U, |U| = N$

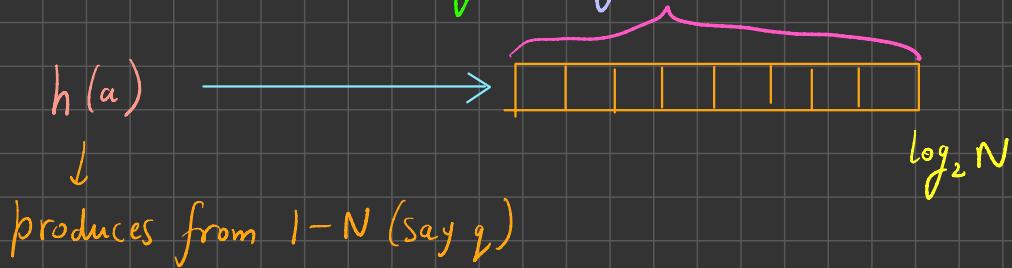
$S = a_0, a_1, \dots, a_n$  stream

$m$  - memory available ( $m \ll n, m \ll N$ )

Given that  $m$  is upto  $\log N$ , we want to approximate distinct elements' count in  $S$ .

## Flajolet-Martin algorithm:

Uses a sketch → bit string  $- \log_2 N$



produces from 1-N (say  $g$ )

Now convert  $g$  to bit string and let  $r(a)$  be the number of trailing zeros in this bit string of  $g$ .

$$\begin{aligned} a_0 &\rightarrow r(a_0) \\ a_1 &\rightarrow r(a_1) \\ a_m &\rightarrow r(a_m) \end{aligned} \quad \left. \right\} R \text{ (max. } r(a_i))$$

Then,

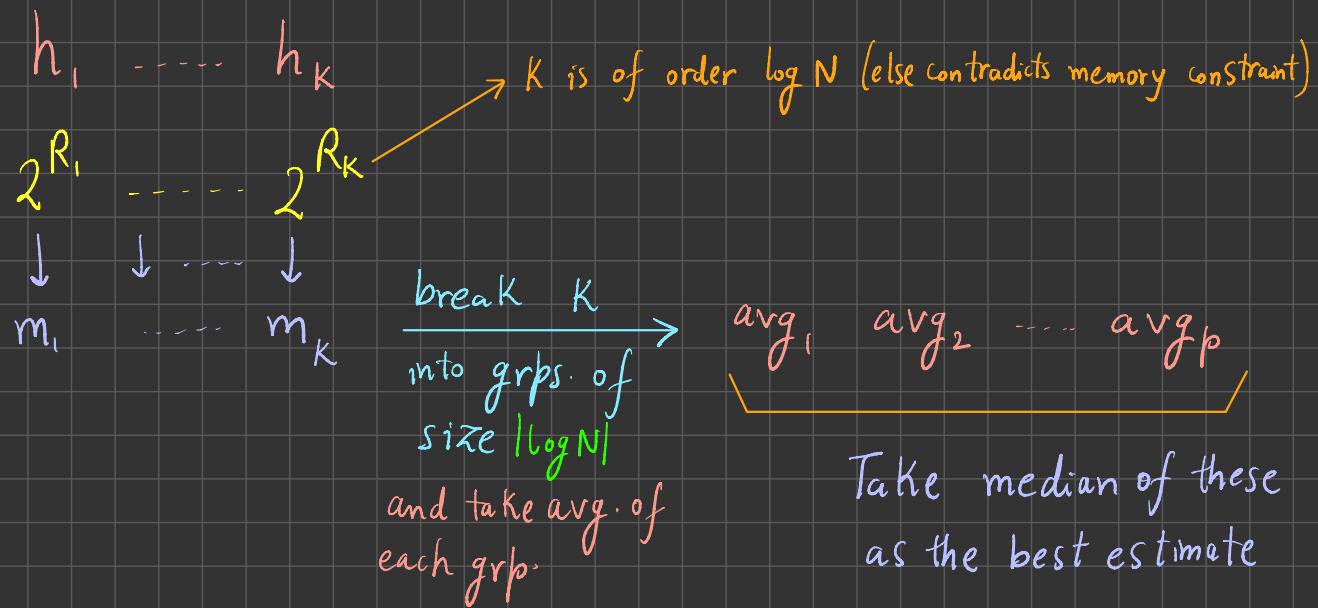
$$\text{approx. count distinct} = 2^R$$

↳ # trailing zeros in bit vector of  $h(a_i)$

more precisely

$$c - \varepsilon \leq \mathbb{E}(2^R) \leq c + \varepsilon$$

Actual count



## Moment Query:

$S = a_0 \dots a_n \quad a_i \in |U|_N \rightarrow U$  is a set with  $|U|=N$

$\downarrow$   
 $v_1, v_2, \dots, v_N$   $\rightarrow$  distinct elements in  $N$ .

$$m_d(S) := \frac{1}{n} \sum_{i=1}^N \left( 1 + c(v_i) + c(v_i)^2 + \dots + c(v_i)^d \right)$$

$c(v_i) = \# v_i \text{ appears in } S$

$$1^{\text{st}} \text{ moment} \rightarrow \frac{\sum c(v_i)}{n}$$

Only  $\log N$  memory available, so we need  
 approx for moment computation.

$$2^{\text{nd}} \text{ moment} \rightarrow \frac{\sum c(v_i)^2}{n}$$

[Surprise number]  
 a number of interest

$X$ -random variable

$$\downarrow \quad x = a_0, a_1, \dots, a_t, \dots, a_n$$

$\rightarrow K$  different values caught at different points of times.

$$S = a_0, \dots, a_t, \dots, a_n$$

$$\xrightarrow{\hspace{1cm}}$$

$$m' \quad n$$

$\hookrightarrow$  frequency of  $x$  over the rest of the stream.

$$C_{\text{approx}}(x) = \frac{n(m')}{n-t}$$

use these for each of  $x_i$  for moment calculations.

General :

$$a \in U, |U| = N$$

$$S = \{a_0, a_1, \dots, a_n\} \longrightarrow g(a)$$

$$|M| = m \quad \left| \begin{array}{l} \text{sketch} \\ \hat{g}(a) \end{array} \right.$$

$$m = O(\log N, \log n) \quad \text{and} \quad \Pr(|\hat{g}(a) - g(a)| \leq \varepsilon) \geq 1 - \delta$$

Frequency query (count)

- additive stream  $\rightarrow$  values once appeared, don't disappear
- add/delete stream

Histogram vector :  $[c_1, \dots, c_N]$

$q(a) = \begin{cases} \text{# of } a & \text{if } a \text{ appears at least } \epsilon n \text{ times in } S \\ 0 & \text{else} \end{cases}$

(Heavy hitters)

↳ So we only want to know the "frequent" terms/heavy hitter in stream

this currently working defn.  
of frequent  
not top K.

Even this is impossible to be answered exactly with memory constraint,  
so an approx  $\hat{q}(a)$  is found.

$$\Pr \left( |\hat{c}(a) - c(a)| \leq \epsilon \right) \geq 1 - \delta$$

Problem: Given a stream  $S$ ,  $|S|=m$ ,  $a \in U$

$$\left. \begin{array}{l} \text{if } c(a) \geq m/2, \text{ then } q(a) = c(a) \\ q(a) = 0 \end{array} \right\}$$

Majority algorithm → works only on additive stream.

Break down of algo:

Initialization data str.

Maintainance / process updates

Query  $q(a)$

## Frequency query (Count) :

Misra-Girres } sketches  
Freq. Count } are linear  
i.e. union of streams  
lead to addition of  
sketches.

Initialization: Counter = 0, value =  $a_0$  (in-hand)

Processing: (i)  $a_i : a_i == \text{value}$ , counter ++

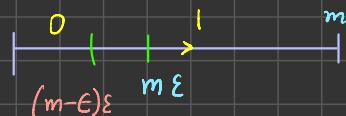
(ii)  $a_i \neq \text{value}$ , counter --

(iii) if counter == 0, value =  $a_i$

Sketch stores  $a_i$ , counter  $\xrightarrow{\text{max size is } m}$   
 $\Rightarrow$  requires  $\log_m$  size.

Problem: if  $c(a) \geq m\epsilon$  then  $g(a) = 1$

$c(a) \leq (m-\epsilon)\epsilon$  then  $g(a) = 0$   
 another epsilon



if  $g(a)$  is 1, find  $\hat{c}(a)$ , such that  $\Pr[|\hat{c}(a) - c(a)| \leq \epsilon'] \geq 1 - \delta$

Misra-Girres algorithm:

$C_1$	$V_1$	$C_2$	$V_2$	$C_3$	$V_2$	$C_4$	$V_4$	$C_K$	$V_K$

Initialization: Counter<sub>1</sub> = 0, Value<sub>1</sub>; Counter<sub>2</sub> = 0, Value<sub>2</sub>; ... . Counter<sub>K</sub> = 0, Value<sub>K</sub>.

Processing: for  $a_i$ ,

if any Counter<sub>p</sub> == 0, update Value<sub>p</sub> =  $a_i$ , and set the corresponding counter to 1.  
 $\xrightarrow{a_{K,p} \text{ already exists}}$

if  $a_i = V_p$ , then increment  
 else decrement all the counters  $\xrightarrow{\text{In space saving algo., we decrement only the one with smallest counter.}}$   
 $\hookrightarrow$  leads to overestimate of count.

Query: if for  $a_i$ , the counter ≠ 0, then  $g(a_i) = 1$

and corresponding counter value is an estimate of frequency

## Count Min Sketch

w, d are independent of stream size.

→ for each stream value, w updates happen

Rest from Slides

C.T. :

Basic Streaming models

Bloom filter

Count distinct

Flajolet Martin

Misra - Gries

Freq. count

Count min

Ref: Ullman / mmmls      } for full proofs  
Graham Cormode lecs.      }  
                                ↳ Sketches, Streaming and  
                                Big data