

Sampling And Streaming

↳ Input is too large to store/access

Plan

- 1) Sampling
- 2) Streaming Model
- 3) Reservoir Sampling
- 4) Distinct Elements

Sampling:

Poll: "Do you have daylight savings?"

YES = 1

NO = 0

Goal: Estimate the fraction of population voting YES

Algo: Sample K - people at random

Collect their answers $X_1, \dots, X_K \in \{0, 1\}$

$$\text{Output} = \frac{1}{K} \sum_{i=1}^K X_i = \hat{p} = \begin{array}{l} \text{estimate} \\ \text{of fraction} \\ \text{of YES} \end{array}$$

Pick k large enough,
With prob. $1 - \delta$ (0.999)

$$\text{estimate } \left| \hat{p} - \underset{p}{\text{true value}} \right| \leq \varepsilon \approx 0.001$$

Question

Suppose for 300 people, K samples are needed
to get 0.01 approximate estimate w.p. 0.999.

How many samples for 300 million people?

1) K samples

2) $100 K$ samples

2) $10 K$

3) $1000 K$ samples

[Chernoff Bound]

Suppose $X_1 \dots X_t$ are $i.i.d$ (independent & identically distributed) random variables that take $\{0,1\}$ values

so that each $\Pr[X_i = 1] = p$ $\Pr[X_i = 0] = 1 - p$

$$\Pr \left[\underbrace{\left| \frac{1}{t} \sum_{i=1}^t X_i - p \right|}_{\substack{\text{average of} \\ \text{\& samples}}} \geq \epsilon \right] \leq 2 e^{-2\epsilon^2 t}$$

\downarrow
expectation

\uparrow
exponentially
in decaying
number
of samples

" δ "

To get: $\Pr[\text{deviation} > \epsilon] \leq \delta$ you pick


$$t = \frac{1}{2\epsilon^2} \log_e \left[\frac{1}{2\delta} \right]$$

$X_i =$ answer of i^{th} person $\in \{0,1\}$

$\Pr(X_i = 1) =$ fraction of YES votes in population
 $= 'p'$

Streaming:



At end of day 

- 1) What fraction of cars are red?
- 2) How many cars travelled?
- 3) How many distinct cars travelled?

- Data :
- 1) Too large to store
 - 2) Comes in a stream,
 - 3) No rewind ~

Streaming Model

Input: a stream of ^{numbers} s_1, \dots, s_n, \dots

$$s_i \in \{1, \dots, N\}$$

[1) read the stream only once, from left
to right]

2) don't know how long the stream is.

Goal:

Algorithm running using space
 $\text{poly}(\log N, \log n)$

Sampling from a Stream

INPUT: A stream $s, \dots \in \{1 \dots N\}$

GOAL: Pick one random element from it.

If $n = \text{length of stream}$ is known

1) Pick a random index $i \in \{1 \dots n\}$

2) Output s_i

Reservoir Sampling

"Reservoir" $v \equiv s_1$

for each element s_i

Pick random number
 $q \in \{1 \dots i\}$

$w = p = 1/i$

$q = i$

$q \neq i$

$w = p = 1 - 1/i$

Throw away what's in reservoir v ,
replace with s_i

Ignore s_i

Whenever stream stops, Output v

Claim: At end of i^{th} step, $\forall j \leq i$

$$\Pr(\text{reservoir} = s_j) = 1/i$$

Proof: Induction:

Base Case = 1 : True

Assume claim true for $i-1$ steps,

After i^{th} step, Consider $j = 1, \dots, i-1$

$$\begin{aligned} \Pr(\text{reservoir} = s_j) &= \Pr(\text{reservoir} = s_j \text{ after } i-1 \text{ steps}) \cdot \Pr(s_j \text{ was not thrown out in } i^{\text{th}} \text{ step}) \\ &= \frac{1}{i-1} \quad // \text{ (induction hypothesis)} \cdot \left(1 - \frac{1}{i}\right) = \frac{1}{i} \end{aligned}$$

$$P_r(\text{reservoir} = s_i) = P_r(\text{replace reservoir with } s_i \text{ in the last step})$$

$$= \frac{1}{n}$$

Distinct Elements:

INPUT: A stream $x_1, \dots, x_n \in \{1, \dots, N\}$

GOAL: Estimate the number of distinct elements in stream.

ALGORITHM (IDEAL)

→ At the beginning, pick a random hash function

$$\left[\underline{h: \{1 \dots N\}} \rightarrow \underbrace{[0, 1]}_{\text{real number}} \right] \times$$

→ Compute (as the stream goes by)

$$\underline{\text{Minimum}}(h(s_1), h(s_2) \dots h(s_n))$$

{ by remembering only one number! }

$$\rightarrow \text{Output} = \left[\begin{array}{c} 1 \\ \text{minimum}(h(s_1) \dots h(s_n)) \end{array} \right]$$

INTUITION: s_1, s_2, \dots, s_n

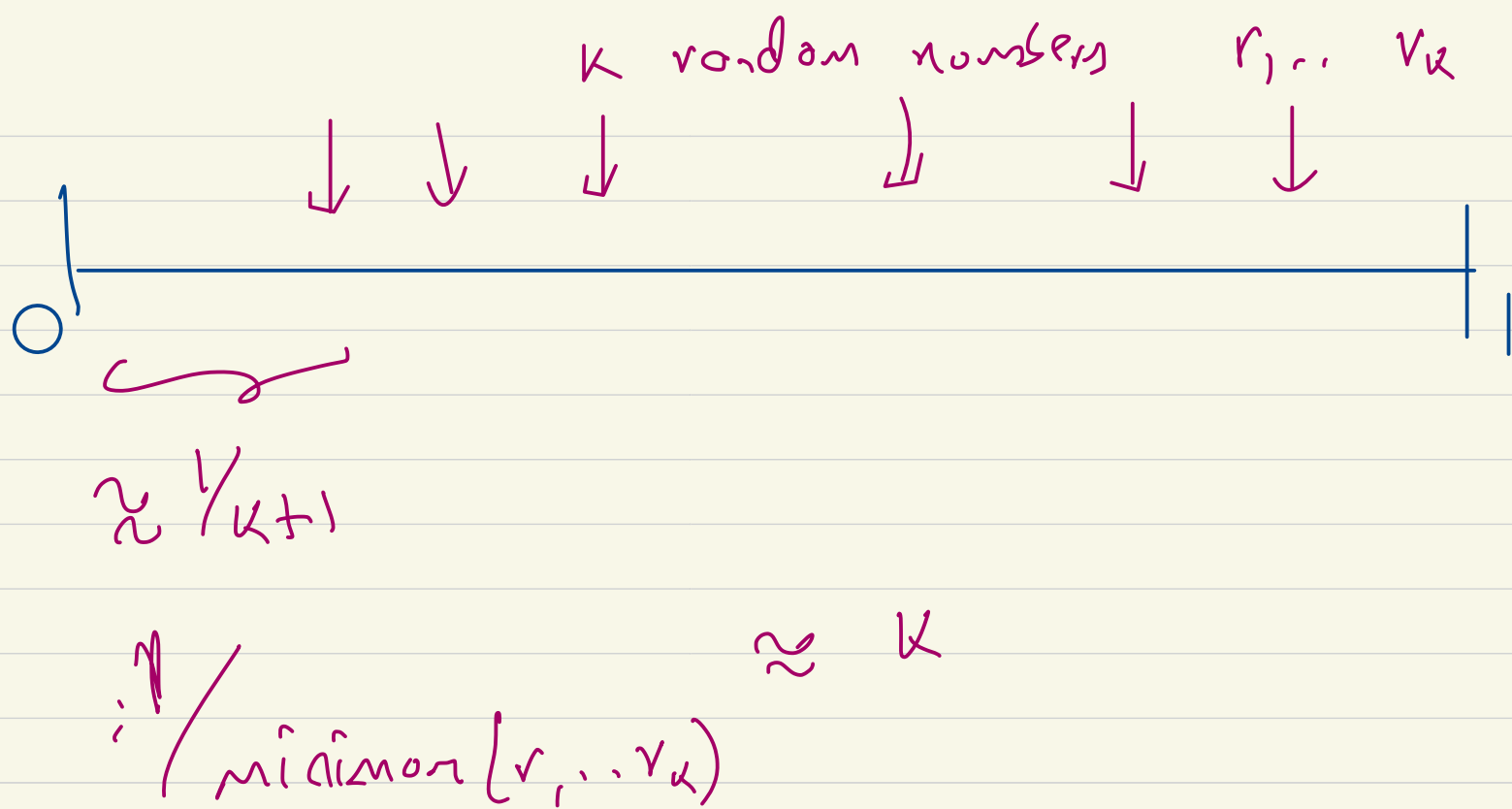
Suppose there are k distinct elements in stream

→ $1, 2, 1, 3, 5, \dots, k$
all numbers

$$\rightarrow \text{minimum}(h(s_1), \dots, h(s_n)) \\ = \text{minimum}(r_1, \dots, r_k)$$

↑
 k different hash values.

[Every hash value is uniformly random i.i.d
in $[0, 1]$]



Lemma:

Pick r_1, \dots, r_k uniformly at random from $[0, 1]$

$$\mathbb{E} \left[\min(r_1, \dots, r_k) \right] = \frac{1}{k+1}$$

Proof: notes.

"Pseudorandom Hash Function"

WANT: For each x_i , $h(x_i) \approx$ uniformly random from $[0,1]$

IMPOSSIBLE:

EXPECT 1: $\begin{matrix} h(1) \\ h(2) \end{matrix} \overset{\text{looks}}{\approx}$ uniformly random in $[0,1]$

"Pairwise independent": $\forall i, j$

$h(i), h(j) \approx$ same distribution as uniformly random pair $[0,1]$