

473 (10/9/18)

Sampling ←
Streaming

Median selection

$B[1 \dots n]$

$R(i)$: i th smallest element
of B .

$Q: R(\frac{n}{2})?$

$$\varepsilon \in (0, \frac{1}{2})$$

$$\Sigma \pm 0.1$$

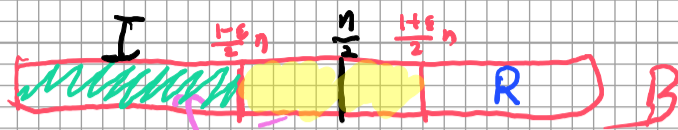
Compute $X \in \mathcal{B}$ s.t.

$$\mathbb{P} \left[R \left((1-\varepsilon)^n \right) \leq X \leq R \left(\frac{1+\varepsilon}{2} n \right) \right] \geq 1 - \delta$$

$\delta \in (0, 1)$ "bad prob".

$$\delta = \frac{1}{n^c}$$

$$K = O \left(\frac{1}{\varepsilon^2} \log \frac{1}{\delta} \right)$$



S sample of size k

$$P = \frac{\left(\frac{1-\epsilon}{2}n\right)}{n}$$

$$= \frac{1-\epsilon}{2}$$

Bad: More than $k/2$ elements of the sample are in I.

$X_i = 1 \Leftrightarrow$ i th sample is
from \bar{I}

$Y = \sum_{i=1}^k X_i$ total number of
samples from \bar{I}

$$P[X_i = 1] = \frac{L-E}{2} = p$$

$Y \geq \frac{k}{2}$

$$\mu = E[Y] = \sum E[X_i] = \frac{L-E}{2} k$$

$$P_1 [Y \geq \frac{k}{2}] \leq P_2 [Y \geq \frac{(1+\epsilon)(1-\epsilon)}{2} k]$$

$$(1+\epsilon)(1-\epsilon) = 1 - \epsilon^2 \leq 1$$

$$P_1 [Y \geq (1+\epsilon)M]$$

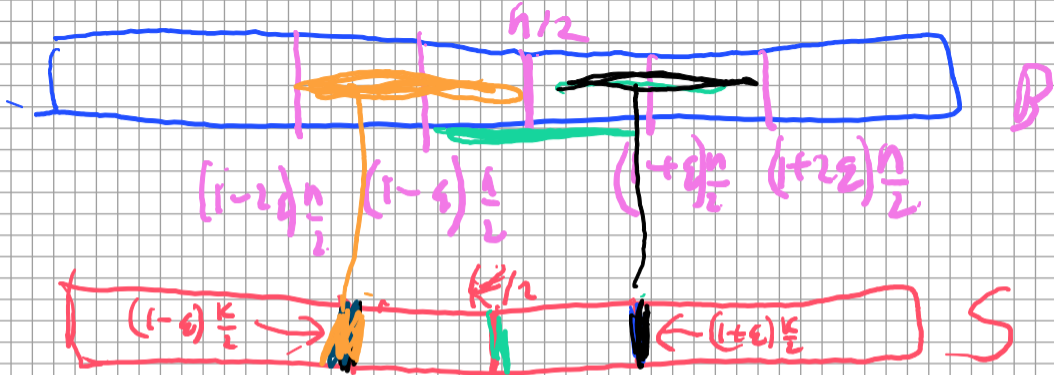
$$P, [Y \geq (1+\epsilon) \mathbb{E}[Y]]$$

Chernoff's inequality

$$\leq \exp\left(-\frac{\epsilon^2 M}{3}\right) \leq \frac{\delta}{2} \quad \square$$

$$M = \frac{1-\epsilon}{2} K \geq \frac{K}{\epsilon}$$

$$K \geq \frac{12}{\epsilon^2} \ln \frac{2}{\delta}$$



$$\kappa = O\left(\frac{1}{\epsilon^2} \log \frac{1}{\sigma}\right)$$

$$P_{\text{err}} \geq \nu \sigma$$



B

$$\left(l - \frac{r}{2} \right)^{\frac{n}{2}} \quad \left(l + \frac{r}{2} \right)^{\frac{n}{2}}$$

$|M| \leq 4\epsilon n \ll n$
 sort M return desired
 element.

$$\varepsilon = \frac{1}{n^{1/4}}$$

$$\sigma = \frac{1}{n^{1/2}}$$

$$k = O\left(\frac{1}{\varepsilon} \log \frac{1}{\sigma}\right) = O(\sqrt{n} \log n)$$

$$O(\sqrt{n} \log n)$$

$$O(\sqrt{n} \log^2 n)$$

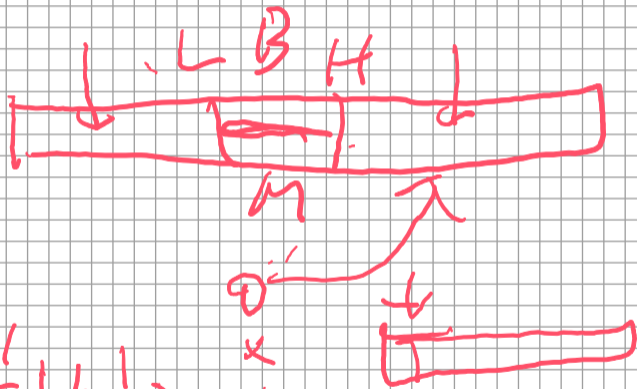
$$|M| = 4 \varepsilon n = O(n^{3/4})$$

$$2n$$
$$O(n^{3/4} \log n)$$

$$2n + O(\sqrt{n} \log^2 n) + O(n^{3/4} \log n)$$

$$2n + o(n)$$

$$LS_n + o(n)$$



$$\frac{1}{2} \cdot 1 + \frac{1}{2} \cdot 2 = \frac{3}{2}$$

Random sampling from a stream

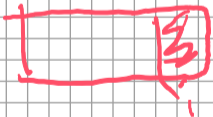
s_1, s_2, \dots

k

Q: Maintain a random sample of size k of the stream
 $O(k)$ space.

→ S: P: $G[0,1]$ $O(1)$
heap with k $O(\log k)$

min-heap to prioritize

$$P_i = \min\left(1, \frac{k}{i}\right)$$


Heavy hitters

$\epsilon_1, \epsilon_2, \dots, \epsilon_n,$

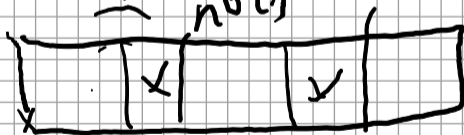
report elements that appear
at least ϵ_n times.

$$k = O\left(\frac{\log n}{\epsilon}\right)$$

random
sample

P [x is a heavy letter
but not in S]

$$\leq \underbrace{\left(1 - \frac{\epsilon}{3}\right)^{k/2}}_{\leq \frac{\epsilon}{3}} = \left(1 - \frac{\epsilon}{3}\right)^{O\left(\frac{\log \frac{1}{\epsilon}}{\frac{\epsilon}{3}}\right)} \leq \frac{1}{n^{O(1)}}$$



ϵ

ϵ_n

S



C



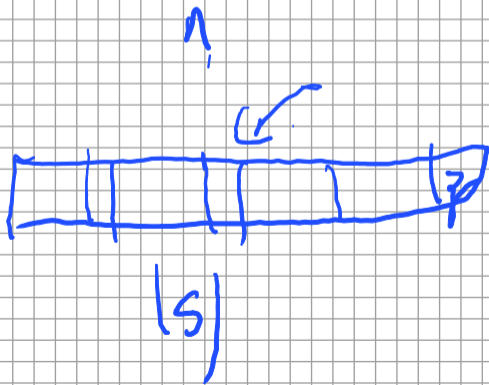
s_i



$$S \left[\begin{array}{c} 2 \\ \sqrt{\frac{1}{\epsilon}} \end{array} \right]$$

total count

Z



alg count

$C(Z)$

$$C(z) = 3$$

$$\text{alg_count}(z) \geq \text{real_count}(z) - \epsilon n$$



