

No class on 26.3.

Reservoir Sampling

Input Stream: e_1, e_2, \dots, e_N

Algorithm Reservoir Sampling (k items without replacement)

Initialize a reservoir array S of size k with the first k items from the stream

Initialize a counter $t \leftarrow k$

for each item x in the stream starting from $(k+1)$ -st item do

$t \leftarrow t + 1$

Draw a random integer j uniformly from $\{1, \dots, t\}$

if $j \leq k$ then

$S[j] \leftarrow x$ // happens with probability $\frac{k}{t}$

return S

Theorem After N items have been processed the reservoir S contains a uniform random sample of k items from the stream without replacement.

Proof: We prove by induction on the stream length $t \geq k$ so far that any item e_i ($i \leq t$) is in the reservoir with probability $\frac{k}{t}$

(Once this claim is established, the theorem follows with $t = N$)

• Base case ($t = k$): The reservoir contains the first k items

$$\text{For } i \leq k, \Pr[e_i \in S] = 1 = \frac{k}{k} = \frac{k}{t} \quad \checkmark$$

• Inductive step ($t \rightarrow t+1$):

Assume that after t items, for any $i \leq t$, $\Pr[e_i \in S \text{ at step } t] = \frac{k}{t}$

Consider $(t+1)$ -st item e_{t+1}

Algorithm includes e_{t+1} in the reservoir with probability $\frac{k}{t+1}$

$$\Pr[e_{t+1} \in S \text{ at step } t+1] = \frac{k}{t+1}$$

Consider item e_i with $i \leq t$

When does e_i that was in the reservoir, stay in the reservoir? ("survives")

• New item e_{t+1} not selected to be in the reservoir (prob. $1 - \frac{k}{t+1}$)

• New item e_{t+1} is selected to be in the reservoir (prob. $\frac{k}{t+1}$),

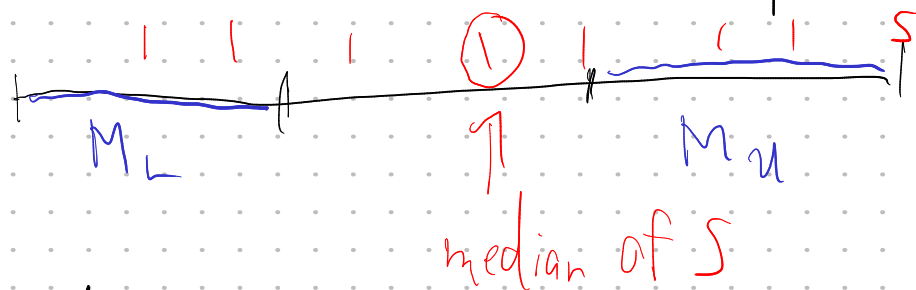
but e_i is not chosen for replacement (with prob. $\frac{k-1}{k}$)

Theorem: By choosing a sample S of size $k = \Theta\left(\frac{1}{\epsilon^2} \cdot \log \frac{1}{\delta}\right)$ in the straight forward algorithm, the median of the sample is an ϵ -approximate median of all items in the stream (i.e., algorithm is successful) with prob. $\geq 1 - \delta$.

Proof:

- M_L : set of smallest $(\frac{1}{2} - \epsilon)N$ elements in the stream
- M_U : ——— " ——— largest ——— " ———

Algorithm fails if median of the sample S is in M_L or M_U



Want to bound $\Pr[\text{median}(S) \in M_L]$ (similar proof then works for M_U)
 $\leq \Pr[\underbrace{\geq \text{half of } S}_{\frac{|S|}{2} \text{ many elements}} \text{ is contained in } M_L] = \Pr[|S_L| > \frac{|S|}{2}]$

• $S_L := S \cap M_L$ $\underline{p} := \Pr[x \in S \text{ is also in } M_L] = \frac{|M_L|}{N} = \frac{(\frac{1}{2} - \epsilon)N}{N} = \frac{1}{2} - \epsilon$

$E[|S_L|] = p \cdot |S| = p \cdot k = k \cdot (\frac{1}{2} - \epsilon)$

$$\begin{aligned}
\Pr[|S_L| > \frac{|S|}{2}] &= \Pr[|S_L| > \frac{k}{2}] \\
&= \Pr[|S_L| - k(\frac{1}{2} - \epsilon) > \frac{k}{2} - k(\frac{1}{2} - \epsilon)] \\
&= \Pr[|S_L| - \mathbb{E}[|S_L|] > \epsilon k] \\
&\leq \frac{1}{e^{\frac{2\epsilon^2 k^2}{k}}} \quad \text{by Hoeffding Bound} \\
&= \frac{1}{e^{2\epsilon^2 k}}
\end{aligned}$$

$$\Pr[\text{Algorithm fails}] \leq 2 \cdot \frac{1}{e^{2\epsilon^2 k}}$$

$$\text{want: } 2 \cdot \frac{1}{e^{2\epsilon^2 k}} \leq \delta$$

$$\Leftrightarrow e^{2\epsilon^2 k} \geq \frac{2}{\delta}$$

$$\Leftrightarrow 2\epsilon^2 k \geq \ln \frac{2}{\delta}$$

$$\Leftrightarrow k \geq \frac{1}{2\epsilon^2} \cdot \ln \frac{2}{\delta}$$

$$\text{setting } k = \lceil \frac{1}{2\epsilon^2} \cdot \ln \frac{2}{\delta} \rceil = \Theta\left(\frac{1}{\epsilon^2} \cdot \log \frac{1}{\delta}\right) \text{ suffices}$$

HW:

Please do some background reading on unbiased estimators

("erwartungstreue Schätzer")