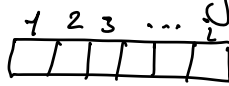


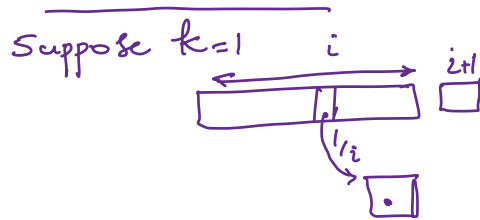
Reservoir Sampling: Unbiased Sampling from Streaming Data.

Given: Streaming data (at any time t , t instances have been observed)



unbiased (Uniform)

Objective: Maintain a Sample of Size k from the Stream



$$u = U_{int}[1, i+1]$$

$$\text{if } (u \leq 1/i)$$

$$\text{Sample} \leftarrow A[i+1]$$

// otherwise, no change

proof (by induction)

$$P_{\text{success}}(A[j]) = \text{Prob. that } A[j] \text{ is the Sample}$$

$$\Rightarrow P_{\text{success}}(A[i+1]) = 1/i+1$$

$$\forall j \leq i \quad P_{\text{success}}(A[j]) = \text{Prob. that it was the Sample at iter } i \text{ AND it Survived the Last Update}$$

$$= 1/i \times (1 - 1/i+1) = 1/i \times \frac{i}{i+1} = \frac{1}{i+1}$$

Extension to $k \geq 1$

for $i = 1$ to k

$$S[i] = A[i]$$

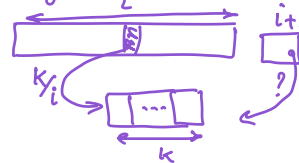
for $i > k$:

$$u = U_{int}[1, i]$$

$$\text{if } (u \leq k/i)$$

$$S[u] = A[i]$$

Proof



$$P_{\text{success}}(A[i+1]) = k/i+1$$

$$\forall j \leq i$$

$$P(A[j] \text{ replaced} \mid \text{in Sample})$$

$$= \frac{1}{i+1}$$

$$\Rightarrow P_{\text{success}}(A[j]) = k/i \times (1 - 1/i+1)$$

$$= k/i+1$$