

①

In this question, we have to present the pseudocode of an algorithm that samples $k \geq 1$ elements uniformly at random from an insert-only stream with unknown length.

Pseudocode:

lst = []

index = 1 // start from 1

while stream[index] exists :

if index $\leq k$:

lst.append(stream[index])

else:

rf = random float from 0 to 1

if rf $< k/\text{index}$:

replace one element from lst with stream[index]

index ++

Proof of Correctness

Base case: index $\leq k$

The list can store k elements. Therefore, with probability of 1, it stores the first k elements from the stream. If the length of stream is less than the k value, we simply return every element in the stream.

Inductive Hypothesis:

We assume that the probability of replacing an element in the list with stream[index] is k/index for every element in the stream after k^{th} element

Induction: ^{index}

The first $\sqrt{\text{index}}$ elements enter the list with probability k/index based on our hypothesis. This means that the next element (index) goes into the list with probability $k/(\text{index}+1)$. If the $(\text{index}+1)$ th element enters the list, we need to take an element out of the list, randomly, which has length of k . Such probability is $1/k$.

Therefore, the probability of the $(\text{index}+1)$ th element entering the list is $1/k (k/(\text{index}+1)) = 1/(\text{index}+1)$.

The complementary, the probability of $(\text{index}+1)$ th element not entering into the list, is therefore

$$1 - 1/(\text{index}+1) = \frac{\text{index}+1-1}{\text{index}+1} = \frac{\text{index}}{\text{index}+1}$$

Finally, the probability that any element stays in the list is $\text{Pr}(\text{not replaced by stream}) \cdot P(\text{entered the list some point})$, which equals

$$\frac{k}{\text{index}} \left(\frac{\text{index}}{\text{index}+1} \right) = \frac{k}{\text{index}+1}, \text{ which matches with our assumption.}$$

Furthermore, as we run through the end of stream, the value goes to

$$\frac{k}{\text{index}} \left(\frac{\text{index}}{\text{index}+1} \right) \left(\frac{\text{index}+1}{\text{index}+2} \right) \cdots \left(\frac{n-1}{n} \right)$$

$$= k/n, \text{ where } n \text{ is the length of the stream.}$$

(2)

Here, we are proving that $\Pr(|Z - Q| \geq \epsilon Q)$ by using Chernoff and Chebyshev bound where $Z = \text{median}_{i \in [t]} \frac{1}{k} \cdot \sum_{j=1}^k x_{ij}$ and

$$t = C_1 \cdot \log\left(\frac{1}{\delta}\right), \quad k = \frac{C_2 \text{Var}[x]}{\epsilon^2 E[x]^2}$$

$$\text{Let } A_i = \frac{1}{k} \cdot \sum_{j=1}^k x_{ij} \quad \text{for } 1 \leq i \leq t$$

$$E[A_i] = Q$$

$$\text{Var}[A_i] = \text{Var}\left(\frac{1}{k} \sum_{j=1}^k x_{ij}\right) = \frac{1}{k^2} \text{Var}\left(\sum_{j=1}^k x_{ij}\right)$$

- Since all x_j are independent RVs,

$$\text{Var}[A_i] = \frac{1}{k^2} \text{Var}\left(\sum_{j=1}^k x_{ij}\right) = \frac{1}{k^2} \sum_{j=1}^k \text{Var}[x] = \frac{1}{k^2} \cdot (k \text{Var}[x]) = \frac{\text{Var}[x]}{k}$$

When we apply Chebyshev inequality, we get

$$\Pr(|A_i - Q| \geq \epsilon Q) \leq \frac{\text{Var}[A_i]}{\epsilon^2 Q^2} = \frac{\text{Var}[x]}{k \epsilon^2 Q^2}$$

$$\text{According to the question, } k = \frac{C_2 \text{Var}[x]}{\epsilon^2 E[x]^2} = \frac{C_2 \text{Var}[x]}{\epsilon^2 Q^2}$$

When we plug in the values,

$$\Pr(|A_i - Q| \geq \epsilon Q) \leq \frac{\text{Var}[x]}{\frac{C_2 \text{Var}[x]}{\epsilon^2 Q^2} \epsilon^2 Q^2} = \frac{1}{C_2}$$

Now, let $C_2 = 4$

Note that Z refers to the median, meaning that Z is greater than half of the values.

$$\hookrightarrow Z > t/2$$

Let $B_i = \begin{cases} 1 & \text{if } |A_i - Q| \geq \epsilon Q \text{ for } 1 \leq i \leq t \\ 0 & \text{otherwise} \end{cases}$

Z , therefore, is $\sum_{i=1}^t B_i$

The expected value of $\sum_{i=1}^t B_i$ is bounded to $t \cdot \frac{1}{2} = t/4$

Since for each B_i where $1 \leq i \leq t$, it is 1 with probability bounded to $1/2 = 1/4$

$$\hookrightarrow np = t \left(\frac{1}{2}\right) = t/4$$

Now,

$\Pr(Z > t/2)$ from the fact that Z is the median

According to Wikipedia and TA,

one version of the Chernoff bound is

$$\Pr(X > (1+\delta')\mu) \leq e^{-\frac{\delta'^2 \mu}{2+\delta'}} \text{ and } \delta' \geq 0$$

Therefore, we need to set $t/2 = (1+\delta')\mu$
 $= (1+\delta')t/4$
 $\delta' = 1$

Substituting the values,

$$\begin{aligned} & \Pr(Z > t/2) \\ &= \Pr(Z > (1+\delta') t/4) < e^{-\frac{(\delta')^2 (t/4)}{2+\delta'}} \quad \text{and } \delta' = 1 \\ &= e^{-\frac{(t/4)}{3}} \\ &= e^{-t/12} = \delta \end{aligned}$$

The question states $t = c_1 \cdot \log(\frac{1}{\delta})$

Plugging in, we get

$$e^{-c_1 \log(\frac{1}{\delta})/12} = \delta$$

In order for this equality to hold, $c_1 = 12$ since $e^{-\log(\frac{1}{\delta})} = \delta$

Therefore, if $c_2 = 4$ and $c_1 = 12$, we proved that $\Pr(|Z - Q| \geq \epsilon Q) \leq \delta$.

$$\textcircled{3} \text{Var}(Z) = E(Z^2) - E(Z)^2$$

From class, we proved that $E(Z) = m+1$

$$\begin{aligned} \text{Var}(Z) &= E(Z^2) - (m+1)^2 \\ &= E(Z^2) - (m^2 + 2m + 1) \end{aligned}$$

Therefore, we need to find $E(Z^2)$

$$\begin{aligned} Z &= 2^{X_m} \\ Z^2 &= (2^{X_m})^2 \\ &= 2^{2X_m} \end{aligned}$$

From the question, we are given that $\text{Var}(Z) = \frac{m(m-1)}{2}$

$$\begin{aligned} E(Z^2) &= \text{Var}(Z) + (m^2 + 2m + 1) \\ &= \frac{m^2 - m}{2} + m^2 + 2m + 1 \end{aligned}$$

$$= \frac{3m^2}{2} + \frac{3m}{2} + 1$$

$$= \frac{3m(m+1)}{2} + 1$$

Therefore we need to prove that $E[2^{2X_m}] = \frac{3m(m+1)}{2} + 1$

(Claim: Define X_m to be the value of the counter after m increments.

$$\text{Then, } E[2^{2X_m}] = \frac{3m(m+1)}{2} + 1$$

Proof: (Induction)

Base case: If $m=0$, $X_m=0$, thus the claim holds

Inductive step:

$$E[2^{2X_m}] = \sum_{j=0}^{\infty} P(X_{m-1}=j) E[2^{2X_m} | X_{m-1}=j]$$

$$= \sum_{j=0}^{\infty} P(X_{m-1}=j) \left[2^{2j} \left(1 - \frac{1}{2^j} \right) + 2^{2(j+1)} \frac{1}{2^j} \right]$$

$$= \sum_{j=0}^{\infty} P(X_{m-1}=j) (2^{2j} - 2^{2j-j} + 2^{2j+2-j})$$

$$= \sum_{j=0}^{\infty} P(X_{m-1}=j) (2^{2j} - 2^j + 2^{j+2})$$

$$= \sum_{j=0}^{\infty} P(X_{m-1}=j) (2^{2j} - 2^j + 4 \cdot 2^j)$$

$$= \sum_{j=0}^{\infty} P(X_{m-1}=j) (2^{2j} + 3 \cdot 2^j)$$

$$= E[2^{2X_{m-1}}] + 3 E[2^{X_{m-1}}]$$

Notice that the next term depends on the first term

$$E[2^{2X_0}] = 1$$

$$E[2^{2X_1}] = 3E[2^{X_0}] + E[2^{2X_0}] = 3 + 1 = 4$$

$$E[2^{2X_2}] = 3E[2^{X_1}] + E[2^{2X_1}] = 3E[2^{X_0}] + 3E[2^{X_0}] + E[2^{2X_0}]$$

$$= 3(2) + 3(1) + 1$$

$$= 6 + 4 = 10$$

↳ Notice a pattern where

$$E[2^{2X_m}] = \sum_{j=1}^m 3j + 1, \text{ which equals } 3 \sum_{j=1}^m j + 1.$$

$$\text{This is } 3 \cdot \frac{(m)(m+1)}{2} + 1.$$

$$\hookrightarrow E[2^{2X_m}] = \frac{3(m)(m+1)}{2} + 1.$$

$$\text{Therefore, } \text{Var}(Z) = \underbrace{\frac{3m(m+1)}{2}}_{E(Z^2)} + 1 - \underbrace{(m^2 + 2m + 1)}_{E(Z)^2} = \frac{m(m-1)}{2}$$

④(a)

Here, we are solving for the pdf of the k -th smallest value among x_1, \dots, x_n for $k=1 \dots n$ where x_1, \dots, x_n are iid uniform random variables and $x_i \in U(0,1)$ for all i .

To find the pdf, we must find the CDF and take the derivative of it.

First, we begin with the smallest hashed value.

Let $V_1 = \min(x_1, \dots, x_n)$

$$E[V_1] = \int_0^1 \Pr(V_1 > t) dt = \int_0^1 \Pr(x_1 > t)^n dt = \int_0^1 (1-t)^n dt = \frac{1}{n+1}$$

In this question, we are concerned in the pdf of k -th smallest hashed value, V_k , where $k \in (1 \dots n)$.

By using the definition of CDF, we are looking for

$\Pr(V_k \leq x)$, which also means $\Pr(\text{at least } k \text{ observations are } \leq x)$

If there are at least k observations, it means that the lower bound is k , which goes to n , the total number of observations.

↳ In addition, for each observation, it follows a binomial distribution with $p = x$.

$$\therefore \Pr(V_k \leq x) = \sum_{l=k}^n \binom{n}{l} x^l (1-x)^{n-l} \rightarrow \text{CDF}$$

If we take the derivative, we get $\frac{d}{dx} \left(\sum_{l=k}^n \binom{n}{l} x^l (1-x)^{n-l} \right)$

$$= \sum_{l=k}^n \binom{n}{l} \cdot \frac{d}{dx} (x^l (1-x)^{n-l})$$

$$= \sum_{l=k}^n \binom{n}{l} (lx^{l-1} \cdot (1-x)^{n-l} - x^l (n-l)(1-x)^{n-l-1})$$

$$= \sum_{l=k}^n \binom{n}{l} (lx^{l-1})(1-x)^{n-l} - \sum_{l=k}^n \binom{n}{l} (x^l)(n-l)(1-x)^{n-l-1}$$

$$= \sum_{l=k}^n \frac{n!}{l!(n-l)!} \cdot l \cdot (x)^{l-1} \cdot (1-x)^{n-l} - \sum_{l=k}^n \frac{n!}{l!(n-l)!} (x^l)(n-l)(1-x)^{n-l-1}$$

↳ when $l=n$, the value of $n-l$ is 0, which makes all of the values to 0 due to multiplication. Therefore, the last term when $n=l$ can be ignored.

$$= \sum_{l=k}^n \frac{n \cdot (n-1)!}{(l-1)!(n-l)!} (x)^{l-1} (1-x)^{n-l} - \sum_{l=k}^{n-1} \frac{n!}{l!(n-l)!} (n-l)(x)^l (1-x)^{n-l-1}$$

$$= \sum_{l=k}^n n \cdot \binom{n-1}{l-1} (x)^{l-1} (1-x)^{n-l} - \sum_{l=k}^{n-1} \frac{n \cdot (n-1)!}{l!(n-l)!} x^l (1-x)^{n-l-1}$$

$$= \sum_{l=k}^n n \binom{n-1}{l-1} (x)^{l-1} (1-x)^{n-l} - \sum_{l=k}^{n-1} n \cdot \binom{n-1}{l} x^l (1-x)^{n-l-1}$$

$$= n \binom{n-1}{k-1} x^{k-1} (1-x)^{(n-1)-(k-1)}$$

$$= \frac{n(n-1)!}{(k-1)!(n-k)!} x^{k-1} (1-x)^{(n-1)-(k-1)}$$

$$= \frac{n!}{(k-1)!(n-k)!} x^{k-1} (1-x)^{(n-1)-(k-1)}$$

(b)

From part a, we found that the pdf of k -th smallest value among x_1, \dots, x_n for $k=1 \dots n$ is $\frac{n!}{(k-1)!(n-k)!} x^{k-1} (1-x)^{(n-1)-(k-1)}$

↳ This is simply the pdf of the beta distribution where

$$\alpha - 1 = k - 1$$

$$\beta - 1 = (n-1) - (k-1)$$

in $B(\alpha, \beta)$

$$\alpha = k$$

$$\beta = n - 1 - k + 1 + 1$$

$$= n - k + 1$$

Therefore, it is the beta distribution $B(k, n-k+1)$.

The expected value of beta distribution $B(\alpha, \beta) = \frac{\alpha}{\alpha + \beta}$

Plugging in, we get $\frac{k}{k + n - k + 1} = \frac{k}{n + 1}$

∴ The expected value of the k -th smallest value among x_1, \dots, x_n for $k=1 \dots n$ is $\frac{k}{n+1}$.