

Can we use the distinct elements (F_0 Estimation) algorithm to sample from distinct elements in a stream with arbitrarily similar probability? — Yes. Run the distinct elements algorithm and at the end of the stream, output a random element from $S \cap U_i$.

Algorithm 2: Sampling distinct elements with arbitrarily similar probability using F_0 estimation

- 1 Set $U_0 = [n]$ and $i = 1$
 - 2 Sample each element of U_{i-1} into U_i with probability $\frac{1}{2}$
 - 3 If $|S \cap U_i| > \frac{2000}{\varepsilon^2} \log n$, then increment $i \leftarrow i + 1$ and repeat from Step 2
 - 4 At the end of the stream, output a random elements from $S \cap U_i$
-

2 Distinct Elements (F_0 Estimation)

Apart from Algorithm 1, another simpler algorithm to estimate the number of distinct elements for insertion-only streams uses hash functions and is as follows:

Algorithm 3:

- 1 Let $h : [n] \rightarrow [0, 1]$ be a random hash function with a real-valued output
 - 2 Initialize $s = 1$
 - 3 For x_1, x_2, \dots, x_m : $s \leftarrow \min(s, h(x_i))$
 - 4 Return $Z = \frac{1}{s} - 1$
-

The intuition behind this algorithm is that the larger the value of N , the smaller we expect s to be.

There are other results that follow from this algorithm:

- It can be shown that $E[s] = \frac{1}{N+1}$ – however this is not the same as $E[Z] = N$ (which is not true!)
- It can be shown that $|s - E[s]| \leq \varepsilon \cdot E[s] \implies (1 - 2\varepsilon)N \leq Z \leq (1 + 4\varepsilon)N$
- It can be shown that $Var[s] \leq \frac{1}{(N+1)^2}$
- It can be shown that by taking the mean of $O\left(\frac{1}{\varepsilon^2}\right)$ independent instances, we get $|s - E[s]| \leq \varepsilon \cdot E[s]$ with probability $\frac{2}{3}$

Note that $(1 - \varepsilon)s \leq E[s] \leq (1 + \varepsilon)s$ implies $(1 - O(\varepsilon)) \cdot N \leq Z \leq (1 + O(\varepsilon)) \cdot N$.

The space complexity of Algorithm 3 is $O(1)$ words. If we run $O\left(\frac{1}{\varepsilon^2}\right)$ independent instances, the space complexity of this algorithm is $O\left(\frac{1}{\varepsilon^2}\right)$ as we only need to store the s value for each of these instances.