

26 Feb 2025

More Streaming Algorithms

Flajolet - Martin Estimator for # Distinct Elements

Let $\mathbb{F}_p = \{0, 1, \dots, p-1\}$

Draw $h: \mathbb{F}_p \rightarrow \mathbb{F}_p$ with random from a 2-universal hash family.

Calculate $Z = \min \{h(x_i)\}$ This calculation requires only $O(1)$ space!

Output the estimate $\frac{p}{Z+1}$.

Last time: If # distinct elements is d ,

$$\Pr(\text{estimate is } > 6d) < \frac{1}{6} \quad [\text{Markov's}]$$

Still need to show -

$$\Pr(\text{estimate} < \frac{d}{6}) < \frac{1}{6} \quad [\text{Chebyshev's}].$$

We defined

$$X_i = \begin{cases} 1 & \text{if } h(x_i) \leq \frac{6d}{d} \\ 0 & \text{o.w.} \end{cases}$$

$$Y = X_1 + \dots + X_d$$

and observed that...

- estimate $< \frac{d}{6}$ when $Y=0$.
- $E[X_i] = \frac{6}{d}$ $\text{Var}(X_i) < E[X_i]$
- $E[Y] = 6 = \frac{6}{d} \cdot d$

$$\bullet \text{Var}(Y) = \sum_{i=1}^d \text{Var}(X_i) < \sum_{i=1}^d E[X_i] = 6.$$

The variance of the sum of pairwise indep rand vars is the sum of their variances.

$$\begin{aligned} \Pr(Y=0) &\leq \Pr(|Y - EY| \geq 6) \\ &= \Pr((Y - EY)^2 \geq 36) \\ &\stackrel{\text{(Markov's)}}{<} \frac{1}{6}. \end{aligned} \quad \left. \begin{array}{l} \text{Chebyshev's} \\ \text{Ineq} \end{array} \right\}$$

$$\dots \text{ So } \Pr\left(\frac{d}{6} \leq \text{estimate} \leq 6d\right) > \frac{2}{3}.$$

To improve accuracy, let $t > 1$.

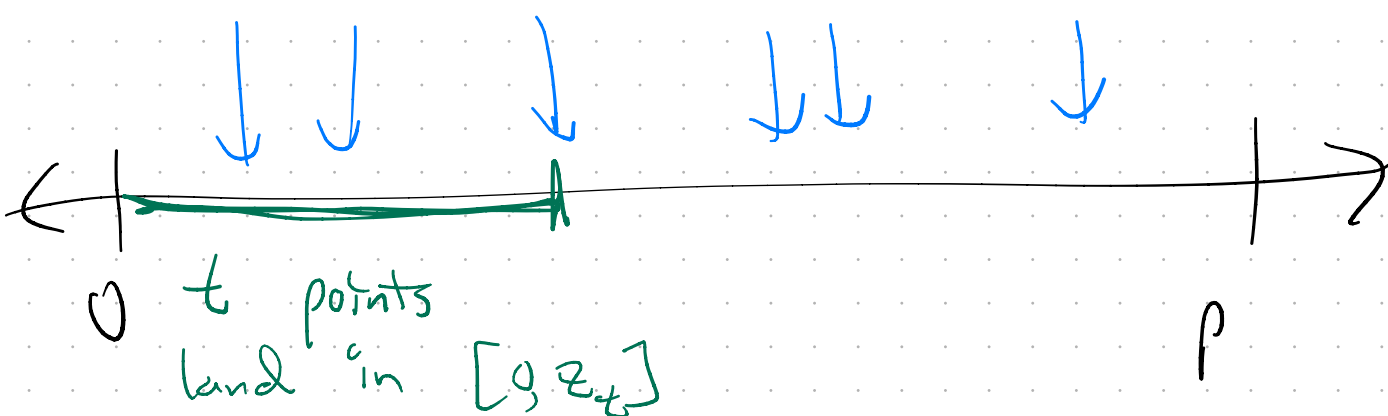
Let $z_1 < z_2 < \dots < z_t$ be the t smallest distinct values among

$\{h(x_1), \dots, h(x_n)\}$.

Estimate is

$$\frac{p \cdot t}{z_t}$$

hash range
estimated density of hash values.



Recap: Assume x_1, \dots, x_d are the d distinct elements.

Again let $d = \#$ distinct elements.

$$X_i = \begin{cases} 1 & \text{if } h(x_i) \leq \frac{(1+\epsilon)t}{d} \\ 0 & \text{otherwise} \end{cases}$$

$$Y = X_1 + \dots + X_d$$

Observe $Y < t$ precisely when

$$Z_t > \frac{(1+\epsilon)t}{d} \Leftrightarrow \frac{t}{Z_t} < \frac{d}{1+\epsilon}$$

$$E[X_i] = \frac{(1+\epsilon)t}{d}$$

$$E[Y] = d \cdot E[X_i] = (1+\epsilon)t$$

$$\text{Var}(Y) = d \cdot \text{Var}(X_i) < d \cdot E[X_i] = (1+\epsilon)t$$

$$\text{Pr}((Y - E[Y])^2 > \epsilon^2 t^2) < \frac{(1+\epsilon)t}{\epsilon^2 t^2}$$

Can make RHS as large as we want by making t large.

To make $\text{Pr}(\text{estimate} < \frac{d}{1+\epsilon}) < \frac{\delta}{2}$,

solve $\frac{1+\epsilon}{\epsilon^2} \cdot \frac{1}{t} < \frac{\delta}{2}$

$$t > \boxed{\frac{2(1+\epsilon)}{\epsilon^2 \delta}}$$

For this value of t , similar
Chebyshev Ineq calculation \Rightarrow

$$\Pr(\text{estimate} > \frac{d}{1-\epsilon}) < \frac{\delta}{2} \checkmark$$

Misra-Gries Algorithm for finding Frequent Elements

- Processes a stream of n
(not necessarily distinct) elements
- Uses $O(k)$ storage
- Reports a list of (at most) k
elements that includes every
element occurring $> \frac{n}{k+1}$ times
in the stream.

... plus false positives.

Maintain list of pairs

(x_i, c_i)



element
seen in
stream

or 1

positive integer

"number of uncancelled occurrences
of x_i in the stream."