

24 Feb 2025

# Streaming Algorithms

## Announcements.

- ① Quiz 3 is graded.
- ② PSet 1 is graded.
- ③ Solution Set 1 on Canvas.
- ④ Quiz 4 on Wed, 2/26, will cover 2/12, 2/19, 2/24.

## Regrade Deadline

Sun 3/2

Sun 3/9

RECAP. A family of hash functions  $\mathcal{H} = \{h: X \rightarrow \mathcal{B}\}$  is 2-universal ("pairwise independent") if

$\forall x \neq y \in X \quad (h(x), h(y)) \in \mathcal{B}^2$  is uniformly distributed.

Ex.  $X = \mathcal{B} = \mathbb{F}_p = \{0, 1, \dots, p-1 \pmod{p}\}$   $p$  prime.  
 $h_{a,b}(x) = ax + b \pmod{p}.$

Why is  $\mathcal{H} = \{h_{a,b} \mid a, b \in \mathbb{F}_p\}$  2-univ.?

For any fixed  $x \neq y$  the pair  $(h(x), h(y))$  can take any value in  $\mathbb{F}_p^2$ ....

because  $h(x)$  can be made to take any value by varying  $b$ ,

and  $h(y) - h(x) = ay + b - (ax + b) = a(y - x)$

can be made to take any value  
by varying  $a$ .

Look at

$$\begin{array}{ll} 0 \cdot (y-x) & \text{mod } p \\ 1 \cdot (y-x) & \text{mod } p \\ \vdots & \end{array}$$

$$(p-1) \cdot (y-x) \text{ mod } p.$$

If any two of these coincide it  
means  $\exists j < k$  in  $\{0, \dots, p-1\}$  s.t.

$$j \cdot (y-x) \equiv k \cdot (y-x) \pmod{p}$$

$$\Rightarrow \underline{(j-k)} \cdot \underline{(y-x)} \text{ is divisible by } p.$$

both factors strictly  
between  $\emptyset$  and  $p$ .

$\Rightarrow$  contradiction. ( $p$  prime)

Inner Product Hashing.

$$X = \mathbb{F}_p^d \quad d > 0$$

$$B = \mathbb{F}_p$$

Def. If  $a = (a_1, \dots, a_d) \in \mathbb{F}_p^d$   
 $x = (x_1, \dots, x_d) \in \mathbb{F}_p^d$

$$\text{let } \langle a, x \rangle = \sum_{i=1}^d a_i x_i \pmod{p}$$

$$h_{\vec{a}, b}(\vec{x}) = \langle \vec{a}, \vec{x} \rangle + b$$

$$h_{\vec{a}, b} : \mathcal{X} \rightarrow \mathcal{B}$$

$$\mathcal{H}_p^d = \left\{ h_{\vec{a}, b} \mid \vec{a} \in \mathbb{F}_p^d, b \in \mathbb{F}_p \right\}$$

Space complexity of representing  
 $h \in \mathcal{H}_p^d$  :  $O(d)$  space to  
 store  $d+1$  values in  $\mathbb{F}_p$

(assume  $p$  is small enough  
 that its binary representation  
 fits in  $O(1)$  memory.)

time complexity of evaluating  
 $h_{\vec{a}, b}$  :  $O(d)$

$d$  mult and  $d$  add in  $\mathbb{F}_p$ .

$\mathbb{H}_p^d$  is  $\mathbb{Z}$ -universal by a variant of the argument already presented.

Consider  $\vec{x} \neq \vec{y} \in \mathbb{F}_p^d$ .

WLOG assume  $x_1 \neq y_1$ .

For any fixed list of coefficients  $a_2, \dots, a_d$  we have  $p^2$  ways to choose  $a_1$  and  $b$ .

By varying  $b$  can make  $h(x)$  take any value in  $\mathbb{F}_p$ .

By varying  $a_1$  can make  $h(y) - h(x)$  take any value in  $\mathbb{F}_p$ .

$$h(y) - h(x) = \langle \vec{a}, \vec{y} - \vec{x} \rangle$$

$$= a_1 \cdot (y_1 - x_1) +$$

ranges over  
all of  $\mathbb{F}_p$   
as  $a_1$  varies.

$$\sum_{i=2}^d a_i (y_i - x_i)$$

constant as  
 $a_2, \dots, a_d, \vec{x}, \vec{y}$   
are fixed.



# Streaming algorithm for distinct elements

$(\epsilon, \delta)$  - PAC algorithm:

for all input sequences / streams

with probability  $\geq 1 - \delta$

the algorithm's answer is correct  
within  $\epsilon$  relative error.

Flajolet - Martin algorithm for  
distinct elements...

[imprecise version,  $\epsilon = 5$ ]

Suppose elements have  $b$ -bit identifiers  
so that  $\{\text{potential elements}\} \subseteq [2^b]$ .

Let  $p$  be a prime  $\geq 2^b$ .

Sample  $h = h_{a,b} \in \mathcal{H}_p^1$  uniformly

at random.

initialize  $Z = p - 1$   
for  $i = 1, \dots, n$ :

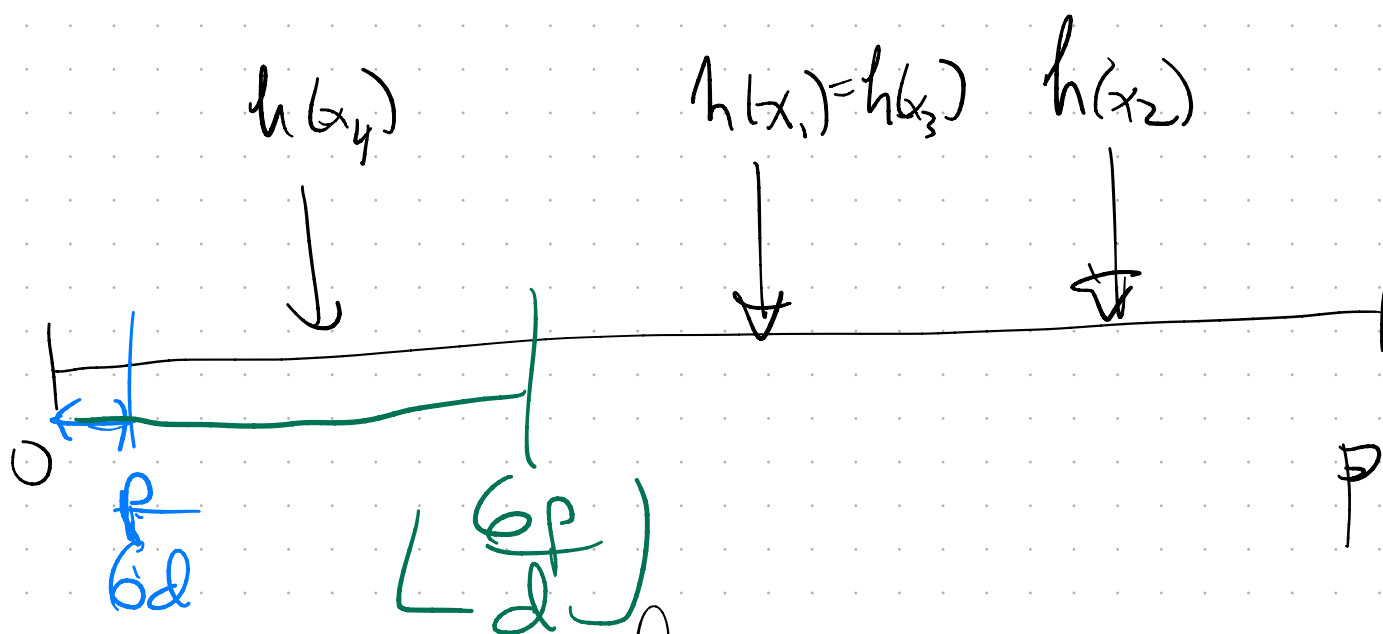
observe  $x_i$  from data stream

compute  $h(x_i)$

update  $Z = \min\{h(x_1), \dots, h(x_i)\}$ .

After observing the stream:

output  $\frac{p}{z+1}$ .



If there are  $d$  distinct values among  $x_1, \dots, x_n$  then you expect  $\{h(x_1), \dots, h(x_n)\}$  to be  $d$  distinct numbers spread uniformly in the interval  $[0, p]$ .

If they were exactly evenly spaced then  $z$  would be at  $\frac{p}{d}$ .

(Maybe  $\frac{p}{d+1}$ )

So  $\frac{p}{z}$  ought <sup>(?)</sup> to be a good estimate of  $d$  or  $d+1$

Analysis. I aim to prove

$$(A) \Pr\left(\frac{p}{z+1} > 6d\right) \leq \frac{1}{6}$$

$$(B) \Pr\left(\frac{f}{Z+1} < \frac{d}{6}\right) \leq \frac{1}{6}$$

$$\text{So then } \Pr\left(\frac{d}{6} \leq \frac{f}{Z+1} \leq 6d\right) \geq \frac{2}{3}$$

...i.e., the algorithm is  $(\epsilon=5, \delta=1/3)$ -PAC

Proof of (A).  $\frac{f}{Z+1} > 6d$

$$\Leftrightarrow Z+1 < \frac{f}{6d}$$

Each individual hash value  $h(x_i)$  is unif distrib over  $\{0, \dots, p-1\}$

$$\Rightarrow \Pr(h(x_i) < \frac{f}{6d}) \leq \frac{1}{6d}$$

There are  $d$  distinct identifiers in the stream

$$\Rightarrow \mathbb{E}\left[\# \text{ identifiers that hash to } < \frac{f}{6d}\right] \leq \left(\frac{1}{6d}\right) \cdot d \leq \frac{1}{6}$$

(linearity of expectation)

$$\Pr(\exists \text{ an identifier hashing to } < \frac{f}{6d}) \leq \frac{1}{6}$$

(Markov's inequality)

Proof of (B), Assume WLOG that  $x_1, \dots, x_d$  are the  $d$  distinct elements in the stream.

Random variables

$$X_{ik} = \begin{cases} 1 & \text{if } h(x_i) \leq k \\ 0 & \text{o.w.} \end{cases}$$

$$Y_k = \sum_{i=1}^d X_{ik}$$

Event (B):  $\frac{p}{Z+1} < \frac{d}{6} \Leftrightarrow Z > \frac{6p}{d} - 1$

Set  $k = \left\lfloor \frac{6p}{d} \right\rfloor$ .

$$Z > \frac{6p}{d} - 1 \Leftrightarrow Y_k = 0.$$

$$\forall i: \mathbb{E}[X_{ik}] \geq \frac{6}{d}.$$

$$\Rightarrow \mathbb{E}[Y_k] \geq 6 \quad (\text{lin of exp.})$$

FACT For a Bernoulli  $\{0,1\}$ -valued

random variable  $X$

$$\begin{aligned}\text{Var}(X) &= \mathbb{E}[X] \cdot (1 - \mathbb{E}[X]) \\ &\leq \mathbb{E}[X].\end{aligned}$$

$$\text{Var}(X_{ik}) \leq \mathbb{E}[X_{ik}]$$

$$\sum_i \text{Var}(X_{ik}) \leq \mathbb{E}\left[\sum X_{ik}\right]$$

pairwise independence  $\Rightarrow$   $\text{Var}(Y_k)$   $\parallel$   $\mathbb{E}[Y_k]$

Lemma. If  $X_1, \dots, X_n$  are pairwise independent and

$$Y = X_1 + \dots + X_n,$$

$$\text{Var}(Y) = \sum_{i=1}^n \text{Var}(X_i)$$

Proof.  $\text{Var}(Y) = E(Y^2) - (E(Y))^2$

$$E(Y) = \sum_{i=1}^n E(X_i)$$

$$(E(Y))^2 = \sum_{i=1}^n \sum_{j=1}^n E(X_i) \cdot E(X_j)$$

$$Y^2 = \sum_{i=1}^n \sum_{j=1}^n X_i X_j$$

$$E(Y^2) = \sum_{i=1}^n \sum_{j=1}^n E[X_i X_j]$$

$$\text{Var}(Y) = \sum_{i=1}^n \sum_{j=1}^n \underbrace{(E[X_i X_j] - E[X_i] \cdot E[X_j])}_{=0 \text{ when } X_i, X_j \text{ independent,}}$$

= 0 when

$X_i, X_j$  independent,

$$= \sum_{i=1}^n E(X_i^2) - (E(X_i))^2$$

$$= \sum_{i=1}^n \text{Var}(X_i)$$

