

1. For every $a \in [\mathbb{P}]^+$ there exists a unique integer $x \in [\mathbb{P}]^+$ such that $ax \bmod p = 1$

Pf: • Since $a \in [\mathbb{P}]^+$, $1 \leq a \leq p-1$, then the greatest common divisor of a and p is 1.

By Bézout's identity, if $\gcd(a, p) = 1$, there exists $x, y \in \mathbb{Z}$ such that $ax + py = 1$.

Reducing module p , we get $ax \equiv 1 \pmod{p}$

• To prove the uniqueness of x , we can use contradiction.

Suppose there is another integer $k \in [\mathbb{P}]^+$ that $ak \equiv 1 \pmod{p}$, then $ax \equiv ak \pmod{p}$

By definition of congruence, $p \mid |ax - ak|$, and 0 is the only non-negative number less than p that is also divisible by p . Thus, $|ax - ak| = 0$

$$\Rightarrow ax = ak$$
$$x = k$$

We thus prove uniqueness that there is no other integer than x for $a \in [\mathbb{P}]^+$ and $ax \bmod p = 1$.

2. To have a hash function that assign inputs uniformly and maximize the number of collision, we can design a hash function that maps all inputs to a single output. For example, let our hash function family as

$$H(x) = 0.$$

This function maps all inputs uniformly to value 0, and all inputs will collide with each other, which results in a maximum number of collisions.

3. Let p be a large prime, $p > |U|$

For any integers $a \in \{1, \dots, p-1\} = [p]^+$, $b \in \{0, 1, \dots, p-1\} = [p]$

define $h_{a,b}(x) = (ax+b) \bmod p \bmod m$.

Let $\mathcal{H} = \{h_{a,b} \mid a \in [p]^+, b \in [p]\}$ be set of all $p(p-1)$ such functions

Prove: \mathcal{H} is 2-universal.

Fix four integers $t_1, t_2, x_1, x_2 \in [p]$ such that $x_1 \neq x_2$,
and $t_1 \neq t_2$.

The linear system $ax_1 + b \equiv t_1 \pmod{p}$ ①

$ax_2 + b \equiv t_2 \pmod{p}$ ②

①-②: $a(x_1 - x_2) \equiv t_1 - t_2 \pmod{p}$

The linear system has a unique solution $a, b \in [p]$ with $a \neq 0$,

where $a = (t_1 - t_2)(x_1 - x_2)^{-1} \bmod p$

$b = (t_2 x_1 - t_1 x_2)(x_1 - x_2)^{-1} \bmod p$.

(notice that b is in a similar expression of a).

Since $x_1 \neq x_2$, a is non-zero if and only if $t_1 \neq t_2$.

By what we prove in problem 1, which implies to here that there is exactly one possible pair of (a, b) that gives us $ax_1 + b = t_1$ and $ax_2 + b = t_2$.

Therefore,

$$\Pr_{a,b} [(ax_1 + b) \bmod p = t_1 \text{ and } (ax_2 + b) \bmod p = t_2] = \frac{1}{p(p-1)}$$

$$\text{and } \Pr_{a,b} [h_{a,b}(x_1) = h_{a,b}(x_2)] = \sum_{i=1}^N \frac{1}{p(p-1)} = \frac{N}{p(p-1)}.$$

where N is the number of ordered pairs $(t_1, t_2) \in [p]^2$ such that $t_1 \neq t_2$, but $t_1 \bmod m = t_2 \bmod m$.

For each fixed $t_1 \in [p]$, there are at most $\lfloor \frac{p}{m} \rfloor$ integers $t_2 \in [p]$ such that $t_1 \neq t_2$ but $t_1 \bmod m = t_2 \bmod m$.

Also, since p is prime, $\lfloor p/m \rfloor \leq (p-1)/m$,

So the number of such ordered pairs: $N \leq \frac{p(p-1)}{m}$.

Therefore, $\Pr(h_{a,b}(x_1) = h_{a,b}(x_2)) \leq \frac{1}{p(p-1)} \cdot \frac{p(p-1)}{m} = \frac{1}{m}$
 which is the condition for 2-universality.

4. Suppose we hash n items into a table of size m .

Let C_{ij} be the indicator variable that equals 1 if and only if $i \neq j$ and $h_{a,b}(i) = h_{a,b}(j)$

Let $C = \sum_{i \neq j} C_{ij}$ be the total number of pairwise collision.

Since hash values are assigned uniformly at random, the probability of a collision is exactly $\frac{1}{m}$.

Linearity of expectation implies:

$$E(C) = \sum_{i \neq j} \Pr(h(i) = h(j)) \leq \binom{n}{2} \frac{2}{m} = \frac{n(n-1)}{m}$$

By Markov's inequality,

$$\Pr(C \geq 1) \leq \frac{E(C)}{1} = 1$$

$$\Pr(C \geq 1) - 1 \leq 1 - 1$$

$$1 - \Pr(C \geq 1) > 0$$

$$\Pr(C = 0) > 0$$

Therefore, there exists a hash function $h \in \mathcal{H}$ that achieves 0 collisions.

5. To prove $\Pr(\text{no slot receives more than } 2 \log n \text{ hashed keys}) \geq 1 - \frac{1}{n}$ is equivalent to prove

$$P(\text{slots receive more than } 2 \log n \text{ hashed keys}) \leq \frac{1}{n}$$

Let X_i be the number of keys hashed to slot i .

$$E(X_i) = n/n = 1$$

Now, considering the prob that slot i contains at least k keys.

There are $\binom{n}{k}$ choices for k keys.

The prob of any particular subset of k keys being hashed in slot i is $\frac{1}{n^k}$, so the union bound ($P(A \cup B) \leq P(A) + P(B)$)

$$\text{implies } P(X_i \geq k) \leq \binom{n}{k} \left(\frac{1}{n}\right)^k \leq \frac{n^k}{k!} \left(\frac{1}{n}\right)^k = \frac{1}{k!}$$

In this case, we consider slots that have at least $2 \log n$ keys, so set $k = 2 \log n$, we have

$$k! \geq 2^k = 2^{2 \log n} = n^2$$

$$\text{which implies } P(X_i \geq 2 \log n) \leq \frac{1}{n^2}$$

This probability bound holds for every slot i .

Thus, by the union bound,

$$\begin{aligned} P(\text{slots receive more than } 2 \log n \text{ hashed keys}) \\ \leq \sum_{i=1}^n P(X_i \geq 2 \log n) \leq n \cdot \frac{1}{n^2} = \frac{1}{n} \end{aligned}$$

Therefore,

$$= 1 - P(\text{slots receive more than } 2 \log n \text{ hashed keys}) \geq 1 - \frac{1}{n}$$

which means

$$\Pr(\text{no slot receives more than } 2 \log n \text{ hashed keys}) \geq 1 - \frac{1}{n}$$

6. To estimate F_2 accurately, we need to compute the sum of products of the frequencies of all 4 elements in the dataset so that we can find a bound for Variance of the output.

Analytically, As we output z^2 for estimation of F_2 , we need to find $E(z^2)$ and $\text{Var}(z^2)$, where we expect variance of z^2 to be tight enough for the estimations.

By linearity of expectation, $\text{Var}(z^2) \leq E(z^4)$, which thus requires 4-wise independence.

To generate hash function for integers, for example, we can design polynomial hash function as

$h(x) = ((ax^3 + bx^2 + cx + d) \bmod p) \bmod m$,
where a, b, c, d are random coefficients chosen uniformly from the set of integers in $[p]^4$. This hash function will return a hash value in the range $[0, m-1]$

- Since storing one integers use 32 bits,

To store a 4 polynomial hash functions with three degrees, we need $4 \times 32 = 128$ bits.

- Alternatively, if we have 64-bits computer where an integer is stored by 64 bits, we need $4 \times 64 = 256$ bits to store the hash function.

7. Theoretical guarantees (proof on blackboard)

Suppose m is the length of the stream. Let the number of buckets B (#cols) be equal to $\lceil \frac{e}{\epsilon} \rceil$

and the number of repetitions (rows) set to $\log(1/\delta)$. Then the estimated frequency \tilde{f}_i

of the true frequency f_i satisfies the following guarantees:

$$1. f_i \leq \tilde{f}_i$$

$$2. \tilde{f}_i \leq f_i + \epsilon m \quad \text{With probability at least } 1-\delta.$$

$$B = \lceil \frac{e}{\epsilon} \rceil$$

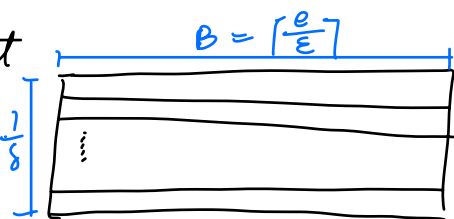
$$r = O(\log(\frac{1}{\delta}))$$

Claim 1: $f_x \leq \tilde{f}_x$ where x is the element

let $\hat{f}_{x,1}, \dots, \hat{f}_{x,r}$ represent the estimated frequency of element x at row r .

Therefore, the estimated frequency

$$\hat{f}_x = \sum_{h(y)=h(x)} f_y = f_x + \sum_{\substack{y \neq x \\ \text{but} \\ \text{hash}(x) = \text{hash}(y)}} f_y \geq f_x$$



$m = \#$ of items in stream

The collision makes the estimated frequency being always larger than the actual frequency

claim 2: if $B = \lceil \frac{e}{\epsilon} \rceil$, $r = \log \frac{1}{\delta}$,

then, $\Pr(\tilde{f}_x \leq f_x + \epsilon m) \geq 1 - \delta$

(\tilde{f}_x is probability not too far from f_x)

proof: Pick item x and define r.v.s $\{Z_1, Z_2, \dots, Z_r\}$ s.t.

$Z_i = \underbrace{C_{i, h(x)}}_{\text{count in row } i} - f_x$ as the over-count in row i due to collisions

let $X_{i,y}$ be the indicator for collision for $i \in \{1, 2, \dots, r\}$, $y \in \{\text{distinct items}\} \setminus \{x\}$

$$X_{i,y} = \begin{cases} 1 & \text{if } h_i(y) = h_i(x) \\ 0 & \text{otherwise} \end{cases}$$

$$\text{Therefore, } Z_i = \sum_{y \neq x} (f_y \cdot X_{i,y})$$

$$E(Z_i) = E\left[\sum_{y \neq x} f_y \cdot X_{i,y}\right] \xrightarrow[\text{expectations}]{\text{linearity of}} \sum_{y \neq x} f_y \cdot E(X_{i,y})$$

$$= \sum_{y \neq x} f_y \cdot \Pr(h_i(y) = h_i(x))$$

Since hash function family H holds universality,

$$\sum_{y \neq x} f_y \cdot \Pr(h_i(y) = h_i(x)) \leq \sum_{y \neq x} f_y \cdot \frac{1}{B} \leq \frac{m}{B}$$

Therefore, expected per-row excess $E(Z_i)$ is at most $\frac{m}{B}$

By Markov's inequality:

$$\Pr(Z_i \geq b \cdot E(Z_i)) \leq \frac{1}{b}, \text{ combine with } E(Z_i) \leq \frac{m}{B}$$

$$\underline{\Pr(Z_i \geq b \cdot \frac{m}{B}) \leq \Pr(Z_i \geq b \cdot E(Z_i)) \leq \frac{1}{b}}$$

$$\hookrightarrow \Pr(Z_i \geq \frac{bm}{B}) \leq \frac{1}{b}$$

Let $b = B\varepsilon$

$$\Pr(Z_i \geq \varepsilon m) \leq \frac{1}{B\varepsilon}$$

$$\text{Since } B = \lceil \frac{e}{\varepsilon} \rceil$$

$$\Pr(Z_i \geq \varepsilon m) \leq \frac{1}{e}$$

$$\text{Therefore, } \Pr(Z_i \geq \varepsilon m) = \Pr(Z_i + f_x \geq f_x + \varepsilon m) \leq \frac{1}{e}$$

If we repeat for r rows and take minimum, and for each row i we perform independently

$$\Pr(\forall 1 \leq i \leq r, Z_i \geq \varepsilon m) \leq \left(\frac{1}{e}\right)^r$$

$$\text{Since } r = \log \frac{1}{\delta}, \left(\frac{1}{e}\right)^r = \left(\frac{1}{e}\right)^{\log \frac{1}{\delta}} = e^{-\log \delta} = \delta$$

$$\text{Therefore, } \Pr(\forall 1 \leq i \leq r, Z_i \geq \varepsilon m) \leq \delta$$

probability of bad estimate is $\leq \delta$

$$\Pr(\exists 1 \leq i \leq r, Z_i < \varepsilon m) \geq 1 - \delta$$

probability of good estimate is $\geq 1 - \delta$

$$\text{so } \tilde{f}_x = \min(C_{1, h_1(x)}, C_{2, h_2(x)}, \dots, C_{r, h_r(x)})$$

$$= \min(f_x + Z_1, f_x + Z_2, \dots, f_x + Z_r) \leq f_x + \varepsilon m \geq 1 - \delta$$