

## Instructions

- The homework is due on **Friday 4/7 at 5pm ET**.
- No extension will be provided, unless for serious documented reasons.
- Start early!
- Study the material taught in class, and feel free to do so in small groups, but the solutions should be a product of your own work.
- This is not a multiple choice homework; reasoning, and mathematical proofs are required before giving your final answer.

## 1 Theory problems [70 pts, 10 each]

In the following let  $p$  be a prime. For any integer  $m$ , define  $[m] = \{0, \dots, m-1\}$  and  $[m]^+ = \{1, \dots, m-1\}$ .

1. Prove that for every  $a \in [p]^+$  there exists a unique integer  $x \in [p]^+$  such that

$$ax \bmod p = 1.$$

2. Answer question 1 on slide 5. Specifically, give a family of hash functions that satisfies the uniformity property but maximizes the number of collisions. Your answer should formally prove why the specific family has the two latter properties.
3. Let  $h_{ab} = (ax + b) \bmod p \bmod m$  where  $a \in [p]^+, b \in [p]$  and  $p$  is a prime such that  $p \geq m$ . Prove that  $\mathcal{H} = \{h_{ab}\}$  is 2-universal.
4. Consider a 2-universal family of hash functions  $\mathcal{H}$  that hash the universe  $U$  to  $[m]$ . Assume you have  $n$  keys  $m > \binom{n}{2}$ . Prove that there exists a hash function  $h \in \mathcal{H}$  that achieves 0 collisions.

$$E[X^2] = E\left[\sum_{i=1}^n f_i^2 Y_i^2 + \sum_{i=1}^n \sum_{j=1, j \neq i}^n f_i f_j Y_i Y_j\right] = E_2$$

*Hint:* Let  $C$  be the RV of number of collisions. Prove that  $Pr_{h \in \mathcal{H}}(C = 0) > 0$ .

5. Suppose we hash  $n$  keys to  $n$  slots. Prove that with probability at least  $1 - \frac{1}{n}$  there is no slot that receives more than  $2 \log n$  hashed keys.
6. Explain why estimating  $F_2$  requires 4-wise independence. Describe how you can generate such as hash function for integers and explain how many bits are needed to store it?
7. In class, we went over the theoretical guarantees (slide 60) of Count-Min sketch when  $B = \lceil \frac{3}{\epsilon} \rceil$  and  $r = O(\log(\frac{1}{\delta}))$  where  $\epsilon, \delta > 0$  are the accuracy and confidence parameters. Your task is the following:

- Write a formal proof of both guarantees 1. and 2. on slide (slide 60). Set the number of buckets  $B = \lceil \frac{e}{\epsilon} \rceil$ .

## **2 Coding [30 points]**

Check the Jupyter notebook on our Git repo.

In the following let  $p$  be a prime. For any integer  $m$ , define  $[m] = \{0, \dots, m-1\}$  and  $[m]^+ = \{1, \dots, m-1\}$ .

1. Prove that for every  $a \in [p]^+$  there exists a unique integer  $x \in [p]^+$  such that

$$ax \bmod p = 1.$$

Proof by contradiction:

Assume there exists more than one integer  $x$  satisfying  $ax \bmod p = 1$ .

Suppose  $ax_1 \bmod p = ax_2 \bmod p$ ,  $x_1, x_2 \in [p]^+$ ,  $x_1 \neq x_2$

$$\Rightarrow a(x_1 - x_2) \bmod p = 0$$

$\because p$  is a prime,  $a \in [1, p-1]$

$\therefore p, a$  are coprime, which means  $(x_1 - x_2)$  can be divided by  $p$ .

$x_1 - x_2 = k \cdot p$ ,  $k$  is an integer.

Yet, we know that  $x \in [1, \dots, p-1]$ , so  $x_1 - x_2 < p$

$\frac{x_1 - x_2}{p} = k < 1$ .  $k$  can't be an integer which is a contradiction.

Thus, for every  $a \in [p]^+$  there exists a unique integer  $x \in [p]^+$

such that  $ax \bmod p = 1$ .



2. Answer question 1 on slide 5. Specifically, give a family of hash functions that satisfies the uniformity property but maximizes the number of collisions. Your answer should formally prove why the specific family has the two latter properties.

Proof:

$$H: \{h_0, h_1, \dots, h_{m-1}\}, \quad h_0(x) = 0, h_1(x) = 1, h_2(x) = 2, \dots, h_{m-1}(x) = m-1$$

- To show it's uniform, we need to show

$$\Pr_{h \in H}(h(x) = i) = \frac{1}{m} \text{ for all } i, x.$$

By definition of  $H$  above,

$$\text{for all } x, \quad h_i(x) = i, \quad \Pr(h(x) = i \mid h(x) = h_i(x)) = 1$$

Since there are  $m$  hash functions in  $H$ ,

$$\Pr_{h \in H}(h(x) = h_i(x)) = \frac{1}{m}$$

$$\Pr(h(x) = h_i(x) \mid h(x) = i) = 1 \quad \begin{array}{l} * \text{ Since there's} \\ \text{no other hash function} \\ \text{hashing } x \text{ to } i \text{ but} \\ h_i(x). \end{array}$$

$$\begin{aligned} \Pr_{h \in H}(h(x) = i) &= \frac{\Pr(h(x) = i \mid h(x) = h_i(x)) \cdot \Pr(h(x) = h_i(x))}{\Pr(h(x) = h_i(x) \mid h(x) = i)} \\ &= \frac{1 \cdot \frac{1}{m}}{1} = \frac{1}{m} \end{aligned}$$

- Collision: items collide if their hash values are equal.  
 $\therefore$  for each item in  $h_i(x)$  their hash values are equal.  
 which is  $h_i(a) = h_i(b) = h_i(c) = \dots = h_i(n) = i$

$\therefore$  The number of collisions = The number of items in the universal set.



3. Let  $h_{ab} = (ax + b) \bmod p \bmod m$  where  $a \in [p]^+, b \in [p]$  and  $p$  is a prime such that  $p \geq m$ . Prove that  $\mathcal{H} = \{h_{ab}\}$  is 2-universal.

Proof:

To prove 2-universal, we need to show that  $\Pr_{h \in \mathcal{H}} (h(x) = h(y)) \leq \frac{1}{m} \quad x \neq y$ .

$a \in \{1, \dots, p-1\}$ ,  $b \in \{0, \dots, p-1\}$ ,  $\Rightarrow a, p$  are coprime,  $b, p$  are coprime.

$r \equiv (ax + b) \bmod p$ ,  $s \equiv (ay + b) \bmod p$ ,  $r \neq s$ .

$\therefore a \neq 0$ ,  $\therefore$  solution  $(a, b)$  is unique.

$$ax \equiv r - b \bmod p$$

$$ay \equiv s - b \bmod p$$

$$ax - ay \equiv (r - b) - (s - b) \bmod p$$

$$a(x - y) \equiv (r - s) \bmod p$$

$$a \equiv (r - s)(x - y)^{-1} \bmod p \quad \because x \neq y$$

$$\left| \begin{array}{l} ax \equiv r - b \bmod p \\ b \equiv r - ax \bmod p \end{array} \right.$$

$$\Rightarrow \begin{array}{l} h_{ab}(x) \equiv r \bmod m \\ h_{ab}(y) \equiv s \bmod m \end{array} \Rightarrow \begin{array}{l} h_{ab}(x) = h_{ab}(y) \\ \text{so } h_{ab}(x) \equiv h_{ab}(y) \bmod m \end{array} \Rightarrow \begin{array}{l} r \equiv s \bmod m, \text{ but} \\ r \neq s \end{array}$$

$$|H| = p \cdot (p-1),$$

$$\Pr_{h \in \mathcal{H}} (h_{ab}(x) = h_{ab}(y)) = \frac{|\overset{\text{collision}}{h_{ab}(x) = h_{ab}(y)}|}{|H|} \leq \frac{\frac{p(p-1)}{m}}{p(p-1)} = \frac{1}{m} \quad \square$$

4. Consider a 2-universal family of hash functions  $\mathcal{H}$  that hash the universe  $U$  to  $[m]$ . Assume you have  $n$  keys  $m > \binom{n}{2}$ . Prove that there exists a hash function  $h \in \mathcal{H}$  that achieves 0 collisions.

Hint: Let  $C$  be the RV of number of collisions. Prove that  $\Pr_{h \in \mathcal{H}}(C = 0) > 0$ .

Proof

$\Pr(h_i(x) = h_i(y), x \neq y) = \frac{1}{m}$ , and for  $n$  keys, there're  $\binom{n}{2}$  pairs of keys that can collide. Thus,  $\Pr(\text{at least a collision}) \leq \binom{n}{2} \cdot \frac{1}{m} < m \cdot \frac{1}{m} = 1$ .

$\Pr(\text{no collision}) = 1 - \Pr(\text{at least a collision}) > 0$

$\therefore$  There exists a hash function that achieves 0 collision.

$\square$

5. Suppose we hash  $n$  keys to  $n$  slots. Prove that with probability at least  $1 - \frac{1}{n}$  there is no slot that receives more than  $2 \log n$  hashed keys.

Proof.

$$\begin{aligned} \Pr(\text{slot } j \text{ has } k \text{ keys}) &= \binom{n}{k} \left(\frac{1}{n}\right)^k \\ &= \frac{n!}{(n-k)! k!} \cdot \frac{1}{n^k} = \frac{n^k}{k! n^k} = \frac{1}{k!} \leq \frac{1}{e^k} \end{aligned}$$

$$\Pr(\text{slot } j \text{ has } > 2 \log n \text{ keys}) = \frac{1}{e^{2 \log n}} = \frac{1}{n^2}$$

$$\begin{aligned} \Pr(\text{no slot } > 2 \log n) &= \Pr(\text{all slots} \leq 2 \log n) = 1 - \Pr(\text{all slots} > 2 \log n) \\ &= 1 - n \cdot \Pr(\text{slot } j > 2 \log n) \leq 1 - \frac{n}{n^2} = 1 - \frac{1}{n} \end{aligned}$$



6. Explain why estimating  $F_2$  requires 4-wise independence. Describe how you can generate such a hash function for integers and explain how many bits are needed to store it?

$$Y_j = h(j) = 1 \text{ or } -1$$

Since  $F_2 = E(X^2) = E\left[\sum_{j=1}^n f_j^2 Y_j^2 + \sum_{i=1}^n \sum_{j=i+1}^n 2 f_i f_j Y_i Y_j\right]$  requires 2-wise independence, in order to make  $E(X^2)$  more accurate, we need to minimize its variance. Thus, 4-wise independence is needed.

$$\begin{aligned} \text{Var}[X^2] &= E[X^4] - (E[X^2])^2 = E[X^4] - F_2^2 \\ E[X^4] &= \sum_{j=1}^n f_j^4 E[Y_j^4] + \binom{4}{2} \sum_{i=1}^n \sum_{j=i+1}^n E[Y_i^2 Y_j^2] f_i f_j \\ &= \sum_{j=1}^n f_j^4 E[Y_j^4] + 3 \left( \left( \sum_{j=1}^n f_j \right)^2 - \sum_{j=1}^n f_j^2 \right) \\ &= F_4 + 3 F_2^2 - 3 F_4 = 3 F_2^2 - 2 F_4 \end{aligned}$$

$$\begin{aligned} \because E[Y_i Y_j] &= \sum Y_i Y_j \Pr(Y_j) = \begin{cases} 1, & Y_i = Y_j \\ 0, & \text{o.w.} \end{cases} \\ \therefore Y_i^2 &= 1, -1^2 = 1 \\ \therefore E[Y_i^2 Y_j^2] &= 1 \\ \therefore E(X^2)^2 &= E[Y_i^2 Y_j^2] \cdot f_i f_j \\ &= \begin{cases} 1 & f_i = f_j \\ 0 & \text{o.w.} \end{cases} \end{aligned}$$

$$\text{Var}[X^2] \leq 3 F_2^2 - 2 F_4 - F_2^2 = 2 F_2^2 - 2 F_4 \leq 2 F_2^2 = 2 \cdot (E[X^2])^2$$

$\Rightarrow 4$  bits.

7. In class, we went over the theoretical guarantees (slide 60) of Count-Min sketch when  $B = \lceil \frac{3}{\epsilon} \rceil$  and  $r = O(\log(\frac{1}{\delta}))$  where  $\epsilon, \delta > 0$  are the accuracy and confidence parameters. Your task is the following:

• Write a formal proof of both guarantees 1. and 2. on slide (slide 60). Set the number of buckets  $B = \lceil \frac{m}{\epsilon} \rceil$ .

Proof:

$$\textcircled{1} \quad f_x \leq \hat{f}_x$$

$$\text{since } \hat{f}_x = f_x + \text{collision},$$

$$\text{we have } f_x \leq \hat{f}_x.$$

$$\textcircled{2} \quad \hat{f}_x \leq f_x + \epsilon m \quad \text{w.p.} \geq 1 - \delta$$

$$E(\text{collision}) \leq \frac{m}{B}$$

$$\because \hat{f}_x = f_x + \text{collision}, \quad \therefore \hat{f}_x \leq f_x + \epsilon m \Rightarrow \Pr(\text{min collision} \geq \epsilon m) \leq \delta$$

$$\Pr(\text{collision}_i \geq \epsilon m) = \frac{E(\text{collision})}{\epsilon m} \leq \frac{\frac{m}{B}}{\epsilon m} = \frac{1}{\epsilon B} \leq \frac{1}{\epsilon} \cdot \frac{\epsilon}{e} = \frac{1}{e} \quad \text{for each row}$$

$$\Pr(\text{min}(\text{collision}_1, \dots, \text{collision}_r) \geq \epsilon m) = \Pr(\text{collision}_1 \geq \epsilon m) \cdot \dots \cdot \Pr(\text{collision}_r \geq \epsilon m) \\ \leq \left(\frac{1}{e}\right)^r = \delta$$

$$\text{Thus, } \hat{f}_x \leq f_x + \epsilon m \quad \text{w.p.} \geq 1 - \delta.$$

