

CS 473: Algorithms

Ruta Mehta

University of Illinois, Urbana-Champaign

Spring 2018

Universal Hashing

Lecture 10

Feb 15, 2018

Most slides are courtesy Prof. Chekuri

Part I

Hash Tables

Dictionary Data Structure

- ① \mathcal{U} : universe of keys with total order: numbers, strings, etc.
- ② Data structure to store a subset $\mathcal{S} \subseteq \mathcal{U}$
- ③ **Operations:**
 - ① **Search/look up:** given $x \in \mathcal{U}$ is $x \in \mathcal{S}$?
 - ② **Insert:** given $x \notin \mathcal{S}$ add x to \mathcal{S} .
 - ③ **Delete:** given $x \in \mathcal{S}$ delete x from \mathcal{S}
- ④ **Static** structure: \mathcal{S} given in advance or changes very infrequently, main operations are lookups.
- ⑤ **Dynamic** structure: \mathcal{S} changes rapidly so inserts and deletes as important as lookups.

Can we do everything in $O(1)$ time?

Hashing and Hash Tables

Hash Table data structure:

- 1 A (hash) table/array T of size m (the table **size**).
- 2 A hash function $h : \mathcal{U} \rightarrow \{0, \dots, m - 1\}$.
- 3 Item $x \in \mathcal{U}$ hashes to slot $h(x)$ in T .

Hashing and Hash Tables

Hash Table data structure:

- 1 A (hash) table/array T of size m (the table **size**).
- 2 A hash function $h : \mathcal{U} \rightarrow \{0, \dots, m - 1\}$.
- 3 Item $x \in \mathcal{U}$ hashes to slot $h(x)$ in T .

Given $S \subseteq \mathcal{U}$. How do we store S and how do we do lookups?

Hashing and Hash Tables

Hash Table data structure:

- 1 A (hash) table/array T of size m (the table **size**).
- 2 A hash function $h : \mathcal{U} \rightarrow \{0, \dots, m - 1\}$.
- 3 Item $x \in \mathcal{U}$ hashes to slot $h(x)$ in T .

Given $S \subseteq \mathcal{U}$. How do we store S and how do we do lookups?

Ideal situation:

- 1 Each element $x \in S$ hashes to a distinct slot in T . Store x in slot $h(x)$
- 2 **Lookup:** Given $y \in \mathcal{U}$ check if $T[h(y)] = y$. $O(1)$ time!

Hashing and Hash Tables

Hash Table data structure:

- 1 A (hash) table/array T of size m (the table **size**).
- 2 A hash function $h : \mathcal{U} \rightarrow \{0, \dots, m - 1\}$.
- 3 Item $x \in \mathcal{U}$ hashes to slot $h(x)$ in T .

Given $S \subseteq \mathcal{U}$. How do we store S and how do we do lookups?

Ideal situation:

- 1 Each element $x \in S$ hashes to a distinct slot in T . Store x in slot $h(x)$
- 2 **Lookup:** Given $y \in \mathcal{U}$ check if $T[h(y)] = y$. $O(1)$ time!

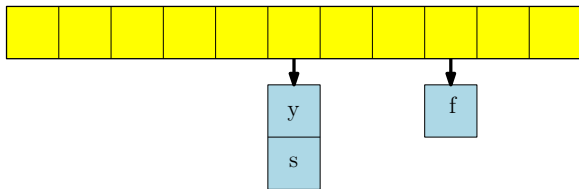
Collisions unavoidable if $|T| < |\mathcal{U}|$. Several techniques to handle them.

Handling Collisions: Chaining

Collision: $h(x) = h(y)$ for some $x \neq y$.

Chaining/Open hashing to handle collisions:

- 1 For each slot i store all items hashed to slot i in a linked list.
 $T[i]$ points to the linked list
- 2 **Lookup:** to find if $y \in \mathcal{U}$ is in T , check the linked list at $T[h(y)]$. Time proportion to size of linked list.



Does hashing give $O(1)$ time per operation for dictionaries?

Hash Functions

Parameters: $N = |\mathcal{U}|$ (very large), $m = |\mathcal{T}|$, $n = |\mathcal{S}|$

Goal: $O(1)$ -time lookup, insertion, deletion.

Single hash function

If $N \geq m^2$, then for any hash function $h : \mathcal{U} \rightarrow \mathcal{T}$ there exists $i < m$ such that at least $N/m \geq m$ elements of \mathcal{U} get hashed to slot i .

Hash Functions

Parameters: $N = |\mathcal{U}|$ (very large), $m = |\mathcal{T}|$, $n = |\mathcal{S}|$

Goal: $O(1)$ -time lookup, insertion, deletion.

Single hash function

If $N \geq m^2$, then for any hash function $h : \mathcal{U} \rightarrow \mathcal{T}$ there exists $i < m$ such that at least $N/m \geq m$ elements of \mathcal{U} get hashed to slot i . Any \mathcal{S} containing all of these is a **very very bad set for h !**

Hash Functions

Parameters: $N = |\mathcal{U}|$ (very large), $m = |\mathcal{T}|$, $n = |\mathcal{S}|$

Goal: $O(1)$ -time lookup, insertion, deletion.

Single hash function

If $N \geq m^2$, then for any hash function $h : \mathcal{U} \rightarrow \mathcal{T}$ there exists $i < m$ such that at least $N/m \geq m$ elements of \mathcal{U} get hashed to slot i . Any \mathcal{S} containing all of these is a **very very bad set for h !**
Such a bad set may lead to $O(m)$ lookup time!

Hash Functions

Parameters: $N = |\mathcal{U}|$ (very large), $m = |\mathcal{T}|$, $n = |\mathcal{S}|$

Goal: $O(1)$ -time lookup, insertion, deletion.

Single hash function

If $N \geq m^2$, then for any hash function $h : \mathcal{U} \rightarrow \mathcal{T}$ there exists $i < m$ such that at least $N/m \geq m$ elements of \mathcal{U} get hashed to slot i . Any \mathcal{S} containing all of these is a **very very bad set for h !**
Such a bad set may lead to $O(m)$ lookup time!

Lesson:

- Consider a family \mathcal{H} of hash functions with *good properties* and choose h uniformly at random.
- Guarantees: small # collisions in expectation for a given \mathcal{S} .
- \mathcal{H} should allow efficient sampling.

Universal Hashing

Question: What are good properties of \mathcal{H} in distributing data?

Universal Hashing

Question: What are good properties of \mathcal{H} in distributing data?

- ① **Uniform:** Consider any element $x \in \mathcal{U}$. Then if $h \in \mathcal{H}$ is picked randomly then x should go into a random slot in \mathcal{T} . In other words $\Pr[h(x) = i] = 1/m$ for every $0 \leq i < m$.

Universal Hashing

Question: What are good properties of \mathcal{H} in distributing data?

- 1 **Uniform:** Consider any element $x \in \mathcal{U}$. Then if $h \in \mathcal{H}$ is picked randomly then x should go into a random slot in \mathcal{T} . In other words $\Pr[h(x) = i] = 1/m$ for every $0 \leq i < m$.
- 2 **Universal:** Consider any two distinct elements $x, y \in \mathcal{U}$. Then if $h \in \mathcal{H}$ is picked randomly then the probability of a collision between x and y should be at most $1/m$. In other words $\Pr[h(x) = h(y)] = 1/m$ (cannot be smaller).

Universal Hashing

Question: What are good properties of \mathcal{H} in distributing data?

- 1 **Uniform:** Consider any element $x \in \mathcal{U}$. Then if $h \in \mathcal{H}$ is picked randomly then x should go into a random slot in T . In other words $\Pr[h(x) = i] = 1/m$ for every $0 \leq i < m$.
- 2 **Universal:** Consider any two distinct elements $x, y \in \mathcal{U}$. Then if $h \in \mathcal{H}$ is picked randomly then the probability of a collision between x and y should be at most $1/m$. In other words $\Pr[h(x) = h(y)] = 1/m$ (cannot be smaller).
- 3 Second property is stronger than the first and the crucial issue.

Definition

A family of hash function \mathcal{H} is (2-)**universal** if for all distinct $x, y \in \mathcal{U}$, $\Pr_{h \sim \mathcal{H}}[h(x) = h(y)] = 1/m$ where m is the table size.

Analyzing Universal Hashing

Question: What is the *expected* time to look up x in T using h assuming chaining used to resolve collisions?

Answer: $O(n/m)$.

Analyzing Universal Hashing

Question: What is the *expected* time to look up x in T using h assuming chaining used to resolve collisions?

Answer: $O(n/m)$.

Comments:

- 1 $O(1)$ expected time also holds for insertion.
- 2 Analysis assumes static set S but holds as long as S is a set formed with at most $O(m)$ insertions and deletions.
- 3 **Worst-case:** look up time can be large! How large?
 $\Omega(\log n / \log \log n)$

Compact Universal Hash Family

Parameters: $N = |\mathcal{U}|$, $m = |\mathcal{T}|$, $n = |\mathcal{S}|$

- 1 Choose a **prime** number $p > N$. Define function
$$h_{a,b}(x) = ((ax + b) \bmod p) \bmod m.$$
- 2 Let $\mathcal{H} = \{h_{a,b} \mid a, b \in \mathbb{Z}_p, a \neq 0\}$ ($\mathbb{Z}_p = \{0, 1, \dots, p - 1\}$).

Compact Universal Hash Family

Parameters: $N = |\mathcal{U}|$, $m = |\mathcal{T}|$, $n = |\mathcal{S}|$

- 1 Choose a **prime** number $p > N$. Define function
$$h_{a,b}(x) = ((ax + b) \bmod p) \bmod m.$$
- 2 Let $\mathcal{H} = \{h_{a,b} \mid a, b \in \mathbb{Z}_p, a \neq 0\}$ ($\mathbb{Z}_p = \{0, 1, \dots, p-1\}$). Note that $|\mathcal{H}| = p(p-1)$.

Compact Universal Hash Family

Parameters: $N = |\mathcal{U}|$, $m = |\mathcal{T}|$, $n = |\mathcal{S}|$

- 1 Choose a **prime** number $p > N$. Define function
$$h_{a,b}(x) = ((ax + b) \bmod p) \bmod m.$$
- 2 Let $\mathcal{H} = \{h_{a,b} \mid a, b \in \mathbb{Z}_p, a \neq 0\}$ ($\mathbb{Z}_p = \{0, 1, \dots, p-1\}$). Note that $|\mathcal{H}| = p(p-1)$.

Theorem

\mathcal{H} is a universal hash family.

Compact Universal Hash Family

Parameters: $N = |\mathcal{U}|$, $m = |\mathcal{T}|$, $n = |\mathcal{S}|$

- 1 Choose a **prime** number $p > N$. Define function $h_{a,b}(x) = ((ax + b) \bmod p) \bmod m$.
- 2 Let $\mathcal{H} = \{h_{a,b} \mid a, b \in \mathbb{Z}_p, a \neq 0\}$ ($\mathbb{Z}_p = \{0, 1, \dots, p-1\}$). Note that $|\mathcal{H}| = p(p-1)$.

Theorem

\mathcal{H} is a universal hash family.

Comments:

- 1 $h_{a,b}$ can be evaluated in $O(1)$ time.
- 2 Easy to store, *i.e.*, just store a, b . Easy to sample.

Some math required...

Lemma (LemmaUnique)

Let p be a prime number, and $\mathbb{Z}_p = \{0, 1, \dots, p - 1\}$.

x : an integer number in \mathbb{Z}_p , $x \neq 0$

\implies There exists a unique $y \in \mathbb{Z}_p$ s.t. $xy = 1 \pmod p$.

In other words: For every element there is a unique inverse.

\implies set $\mathbb{Z}_p = \{0, 1, \dots, p - 1\}$ when working modulo p is a field.

Proof of Lemma Unique

Claim

Let p be a prime number. For any $x, y, z \in \{1, \dots, p-1\}$ s.t. $y \neq z$, we have that $xy \bmod p \neq xz \bmod p$.

Proof of Lemma Unique

Claim

Let p be a prime number. For any $x, y, z \in \{1, \dots, p-1\}$ s.t. $y \neq z$, we have that $xy \bmod p \neq xz \bmod p$.

Proof.

Assume for the sake of contradiction $xy \bmod p = xz \bmod p$.
Then

$$x(y - z) = 0 \bmod p$$

$$\implies p \text{ divides } x(y - z)$$

$$\implies p \text{ divides } x \text{ OR } p \text{ divides } (y - z) \quad (\text{why?})$$

$$\implies y - z = 0 \implies y = z$$

And that is a contradiction. □

Proof of LemmaUnique

Lemma (LemmaUnique)

Let p be a prime number,

x : an integer number in $\{1, \dots, p - 1\}$.

\implies There exists a unique y s.t. $xy = 1 \pmod p$.

Proof.

By the above claim if $xy = 1 \pmod p$ and $xz = 1 \pmod p$ then $y = z$. Hence uniqueness follows.

Proof of LemmaUnique

Lemma (LemmaUnique)

Let p be a prime number,

x : an integer number in $\{1, \dots, p - 1\}$.

\implies There exists a unique y s.t. $xy = 1 \pmod p$.

Proof.

By the above claim if $xy = 1 \pmod p$ and $xz = 1 \pmod p$ then $y = z$. Hence uniqueness follows.

Existence. For any $x \in \{1, \dots, p - 1\}$ we have that $\{x * 1 \pmod p, x * 2 \pmod p, \dots, x * (p - 1) \pmod p\} =$

Proof of LemmaUnique

Lemma (LemmaUnique)

Let p be a prime number,

x : an integer number in $\{1, \dots, p - 1\}$.

\implies There exists a unique y s.t. $xy = 1 \pmod p$.

Proof.

By the above claim if $xy = 1 \pmod p$ and $xz = 1 \pmod p$ then $y = z$. Hence uniqueness follows.

Existence. For any $x \in \{1, \dots, p - 1\}$ we have that $\{x * 1 \pmod p, x * 2 \pmod p, \dots, x * (p - 1) \pmod p\} = \{1, 2, \dots, p - 1\}$.

\implies There exists a number $y \in \{1, \dots, p - 1\}$ such that $xy = 1 \pmod p$. □

Proof of the Theorem: Outline

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

Theorem

$\mathcal{H} = \{h_{a,b} \mid a, b \in \mathbb{Z}_p, a \neq 0\}$ is universal.

Proof.

Fix $x, y \in \mathcal{U}$. Show that $\Pr_{h_{a,b} \sim \mathcal{H}}[h_{a,b}(x) = h_{a,b}(y)] \leq 1/m$.

Note that $|\mathcal{H}| = p(p-1)$.

Proof of the Theorem: Outline

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

Theorem

$\mathcal{H} = \{h_{a,b} \mid a, b \in \mathbb{Z}_p, a \neq 0\}$ is universal.

Proof.

Fix $x, y \in \mathcal{U}$. Show that $\Pr_{h_{a,b} \sim \mathcal{H}}[h_{a,b}(x) = h_{a,b}(y)] \leq 1/m$.

Note that $|\mathcal{H}| = p(p-1)$.

- Let (a, b) (equivalently $h_{a,b}$) be *bad* for x, y if $h_{a,b}(x) = h_{a,b}(y)$.

Proof of the Theorem: Outline

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

Theorem

$\mathcal{H} = \{h_{a,b} \mid a, b \in \mathbb{Z}_p, a \neq 0\}$ is universal.

Proof.

Fix $x, y \in \mathcal{U}$. Show that $\Pr_{h_{a,b} \sim \mathcal{H}}[h_{a,b}(x) = h_{a,b}(y)] \leq 1/m$.

Note that $|\mathcal{H}| = p(p-1)$.

- 1 Let (a, b) (equivalently $h_{a,b}$) be *bad* for x, y if $h_{a,b}(x) = h_{a,b}(y)$. At most howmany bad h is ok?

Proof of the Theorem: Outline

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

Theorem

$\mathcal{H} = \{h_{a,b} \mid a, b \in \mathbb{Z}_p, a \neq 0\}$ is universal.

Proof.

Fix $x, y \in \mathcal{U}$. Show that $\Pr_{h_{a,b} \sim \mathcal{H}}[h_{a,b}(x) = h_{a,b}(y)] \leq 1/m$.

Note that $|\mathcal{H}| = p(p-1)$.

- 1 Let (a, b) (equivalently $h_{a,b}$) be *bad* for x, y if $h_{a,b}(x) = h_{a,b}(y)$. At most howmany bad h is ok?
- 2 **Claim:** Number of bad (a, b) is at most $p(p-1)/m$.

Proof of the Theorem: Outline

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

Theorem

$\mathcal{H} = \{h_{a,b} \mid a, b \in \mathbb{Z}_p, a \neq 0\}$ is universal.

Proof.

Fix $x, y \in \mathcal{U}$. Show that $\Pr_{h_{a,b} \sim \mathcal{H}}[h_{a,b}(x) = h_{a,b}(y)] \leq 1/m$.

Note that $|\mathcal{H}| = p(p-1)$.

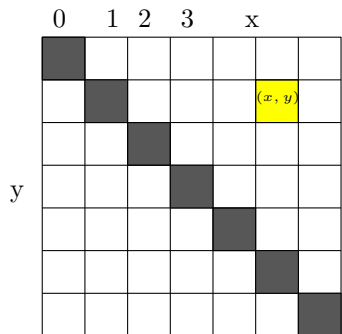
- 1 Let (a, b) (equivalently $h_{a,b}$) be *bad* for x, y if $h_{a,b}(x) = h_{a,b}(y)$. At most howmany bad h is ok?
- 2 **Claim:** Number of bad (a, b) is at most $p(p-1)/m$.
- 3 Total number of hash functions is $p(p-1)$ and hence probability of a collision is $\leq 1/m$. □

Intuition for the Claim

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

First map $x \neq y$ to $r = g_{a,b}(x)$ and $s = g_{a,b}(y)$.

Lemma Unique proof $\implies r \neq s$

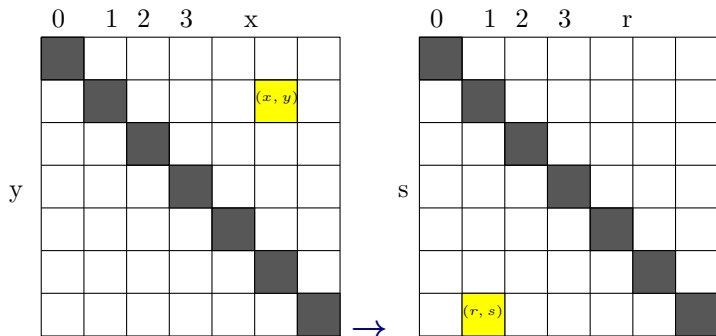


Intuition for the Claim

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

First map $x \neq y$ to $r = g_{a,b}(x)$ and $s = g_{a,b}(y)$.

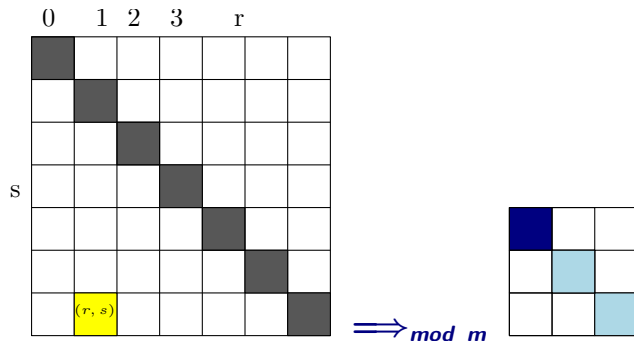
LemmaUnique proof $\implies r \neq s$



As (a, b) varies, (r, s) takes all possible $p(p - 1)$ values. Since (a, b) is picked u.a.r., every value of (r, s) has equal probability.

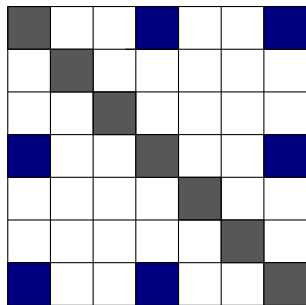
Intuition for the Claim

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

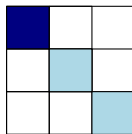


Intuition for the Claim

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



$\implies \bmod m$



Intuition for the Claim

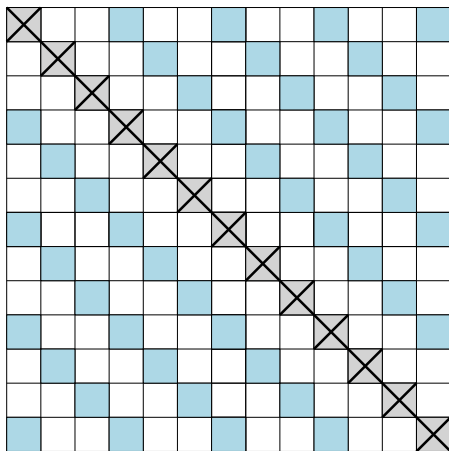
For a fixed $a \in \{0, \dots, m - 1\}$ what is an upper bound on the size of set $\{s \in \{0, \dots, (p - 1)\} \mid a = s \pmod m\}$?

- (A) m .
- (B) m^2 .
- (C) p .
- (D) p/m .
- (E) Many. At least two.

Intuition for the Claim

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

- 1 First part of mapping maps (x, y) to a random location $(g_{a,b}(x), g_{a,b}(y))$ in the “matrix”.
- 2 $(g_{a,b}(x), g_{a,b}(y))$ is not on main diagonal.
- 3 All blue locations are “bad” – map by $\bmod m$ to a location of collision.
- 4 But... at most $1/m$ fraction of allowable locations in the matrix are bad.



We need

to show at most $1/m$ fraction of bad $h_{a,b}$

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

2 lemmas ...

Fix $x \neq y \in \mathbb{Z}_p$, and let $r = (ax + b) \bmod p$ and $s = (ay + b) \bmod p$.

We need

to show at most $1/m$ fraction of bad $h_{a,b}$

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

2 lemmas ...

Fix $x \neq y \in \mathbb{Z}_p$, and let $r = (ax + b) \bmod p$ and $s = (ay + b) \bmod p$.

- 1 1-to-1 correspondence between $p(p - 1)$ pairs of (a, b) (equivalently $h_{a,b}$) and $p(p - 1)$ pairs of (r, s) .

We need

to show at most $1/m$ fraction of bad $h_{a,b}$

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

2 lemmas ...

Fix $x \neq y \in \mathbb{Z}_p$, and let $r = (ax + b) \bmod p$ and $s = (ay + b) \bmod p$.

- 1-to-1 correspondence between $p(p-1)$ pairs of (a, b) (equivalently $h_{a,b}$) and $p(p-1)$ pairs of (r, s) .
- Out of all possible $p(p-1)$ pairs of (r, s) , at most $p(p-1)/m$ fraction satisfies $r \bmod m = s \bmod m$.

Some Lemmas

Lemma

If $x \neq y$ then for any $a, b \in \mathbb{Z}_p$ such that $a \neq 0$, we have
$$ax + b \pmod p \neq ay + b \pmod p.$$

Some Lemmas

Lemma

If $x \neq y$ then for any $a, b \in \mathbb{Z}_p$ such that $a \neq 0$, we have
 $ax + b \pmod p \neq ay + b \pmod p$.

Proof.

Suppose not

$$ax + b \pmod p = ay + b \pmod p \Rightarrow a(x - y) \pmod p = 0$$

Some Lemmas

Lemma

If $x \neq y$ then for any $a, b \in \mathbb{Z}_p$ such that $a \neq 0$, we have
$$ax + b \pmod p \neq ay + b \pmod p.$$

Proof.

Suppose not

$$ax + b \pmod p = ay + b \pmod p \Rightarrow a(x - y) \pmod p = 0$$

Since p is a prime, p divides either a or $(x - y)$.

Some Lemmas

Lemma

If $x \neq y$ then for any $a, b \in \mathbb{Z}_p$ such that $a \neq 0$, we have
 $ax + b \pmod p \neq ay + b \pmod p$.

Proof.

Suppose not

$$ax + b \pmod p = ay + b \pmod p \Rightarrow a(x - y) \pmod p = 0$$

Since p is a prime, p divides either a or $(x - y)$. But $a < p$ and $(x - y) < p$, and hence $a = 0$ or $(x - y) = 0$. Contradiction! \square

Some Lemmas

Lemma

If $x \neq y$ then for each (r, s) such that $r \neq s$ and $0 \leq r, s \leq p - 1$ there is exactly **one** a, b such that
 $ax + b \pmod p = r$ and $ay + b \pmod p = s$.

Proof.

Solve the two equations:

$$ax + b = r \pmod p \quad \text{and} \quad ay + b = s \pmod p$$

Some Lemmas

Lemma

If $x \neq y$ then for each (r, s) such that $r \neq s$ and $0 \leq r, s \leq p - 1$ there is exactly **one** a, b such that $ax + b \pmod p = r$ and $ay + b \pmod p = s$.

Proof.

Solve the two equations:

$$ax + b = r \pmod p \quad \text{and} \quad ay + b = s \pmod p$$

We get $a = \frac{r-s}{x-y} \pmod p$ and $b = r - ax \pmod p$. □

One-to-one correspondence between (a, b) and (r, s)

Understanding the hashing

Once we fix a and b , and we are given a value x , we compute the hash value of x in two stages:

- 1 **Compute:** $r \leftarrow (ax + b) \bmod p$.
- 2 **Fold:** $r' \leftarrow r \bmod m$

Understanding the hashing

Once we fix a and b , and we are given a value x , we compute the hash value of x in two stages:

- 1 **Compute:** $r \leftarrow (ax + b) \bmod p$.
- 2 **Fold:** $r' \leftarrow r \bmod m$

Collision...

Given two distinct values x and y they might collide only because of folding.

Lemma

not equal pairs (r, s) of $\mathbb{Z}_p \times \mathbb{Z}_p$ that are folded to the same number is $p(p - 1)/m$.

Folding numbers

Lemma

pairs $(r, s) \in \mathbb{Z}_p \times \mathbb{Z}_p$ such that $r \neq s$ and $r \bmod m = s \bmod m$ (folded to the same number) is $p(p-1)/m$.

Proof.

Consider a pair $(r, s) \in \{0, 1, \dots, p-1\}^2$ s.t. $r \neq s$. Fix r :

① $a = r \bmod m$.

Folding numbers

Lemma

pairs $(r, s) \in \mathbb{Z}_p \times \mathbb{Z}_p$ such that $r \neq s$ and $r \bmod m = s \bmod m$ (folded to the same number) is $p(p-1)/m$.

Proof.

Consider a pair $(r, s) \in \{0, 1, \dots, p-1\}^2$ s.t. $r \neq s$. Fix r :

- 1 $a = r \bmod m$.
- 2 There are $\lceil p/m \rceil$ values of s that fold into a . That is

$$r \bmod m = s \bmod m.$$

- 3 One of them is when $r = s$.
- 4 \implies # of colliding pairs

Folding numbers

Lemma

pairs $(r, s) \in \mathbb{Z}_p \times \mathbb{Z}_p$ such that $r \neq s$ and $r \bmod m = s \bmod m$ (folded to the same number) is $p(p-1)/m$.

Proof.

Consider a pair $(r, s) \in \{0, 1, \dots, p-1\}^2$ s.t. $r \neq s$. Fix r :

- 1 $a = r \bmod m$.
- 2 There are $\lceil p/m \rceil$ values of s that fold into a . That is

$$r \bmod m = s \bmod m.$$

- 3 One of them is when $r = s$.
- 4 \implies # of colliding pairs $(\lceil p/m \rceil - 1)p \leq (p-1)p/m$



Proof of Claim

of bad pairs is $p(p-1)/m$

Proof.

Let $a, b \in \mathbb{Z}_p$ such that $a \neq 0$ and $h_{a,b}(x) = h_{a,b}(y)$.

- ① Let $r = ax + b \pmod p$ and $s = ay + b \pmod p$.
- ② Collision if and only if $r \pmod m = s \pmod m$.
- ③ (Folding error): Number of pairs (r, s) such that $r \neq s$ and $0 \leq r, s \leq p-1$ and $r \pmod m = s \pmod m$ is $p(p-1)/m$.
- ④ From previous lemma there is one-to-one correspondence between (a, b) and (r, s) . Hence total number of bad (a, b) pairs is $p(p-1)/m$.



Proof of Claim

of bad pairs is $p(p-1)/m$

Proof.

Let $a, b \in \mathbb{Z}_p$ such that $a \neq 0$ and $h_{a,b}(x) = h_{a,b}(y)$.

- ① Let $r = ax + b \pmod p$ and $s = ay + b \pmod p$.
- ② Collision if and only if $r \pmod m = s \pmod m$.
- ③ (Folding error): Number of pairs (r, s) such that $r \neq s$ and $0 \leq r, s \leq p-1$ and $r \pmod m = s \pmod m$ is $p(p-1)/m$.
- ④ From previous lemma there is one-to-one correspondence between (a, b) and (r, s) . Hence total number of bad (a, b) pairs is $p(p-1)/m$.



Prob of x and y to collide: $\frac{\# \text{ bad } (a, b) \text{ pairs}}{\#(a, b) \text{ pairs}} = \frac{p(p-1)/m}{p(p-1)} = \frac{1}{m}$.

Look up Time

Say $|S| = |T| = m$.

For $0 \leq i \leq m - 1$, $\ell(i)$: list of elements hashed to slot i in T .

Expected look up time

Since for $x \neq y$, $\Pr[h_{a,b}(x) = h_{a,b}(y)] = 1/m$, we get
 $\mathbf{E}[|\ell(i)|] = |S|/m \leq 1$.

Look up Time

Say $|S| = |T| = m$.

For $0 \leq i \leq m - 1$, $\ell(i)$: list of elements hashed to slot i in T .

Expected look up time

Since for $x \neq y$, $\Pr[h_{a,b}(x) = h_{a,b}(y)] = 1/m$, we get
 $\mathbf{E}[|\ell(i)|] = |S|/m \leq 1$.

Expected worst case look up time

Like in *Balls & Bins*, $\mathbf{E}\left[\max_{i=0}^{m-1} |\ell(i)|\right] \geq O(\ln n / \ln \ln n)$.

Look up Time

Say $|S| = |T| = m$.

For $0 \leq i \leq m - 1$, $\ell(i)$: list of elements hashed to slot i in T .

Expected look up time

Since for $x \neq y$, $\Pr[h_{a,b}(x) = h_{a,b}(y)] = 1/m$, we get
 $\mathbf{E}[|\ell(i)|] = |S|/m \leq 1$.

Expected worst case look up time

Like in *Balls & Bins*, $\mathbf{E}\left[\max_{i=0}^{m-1} |\ell(i)|\right] \geq O(\ln n / \ln \ln n)$.

What if $|T| = m^2$ (# Bins is m^2)

Claim: If $|T| = m^2$, then $\mathbf{E}\left[\max_{i=0}^{m-1} |\ell(i)|\right] = O(1)$.

Perfect Hashing

Two levels of hash tables

Question: Can we make look up time $O(1)$ in worst case?

Perfect Hashing for Static Data

- Do hashing once.
- If $Y_i = |\ell(i)| > 10$ then hash elements of $\ell(i)$ to a table of size Y_i^2 .

Perfect Hashing

Two levels of hash tables

Question: Can we make look up time $O(1)$ in worst case?

Perfect Hashing for Static Data

- Do hashing once.
- If $Y_i = |\ell(i)| > 10$ then hash elements of $\ell(i)$ to a table of size Y_i^2 .

Lemma (Look-up)

Expected worst case look up time is $O(1)$.

Perfect Hashing

Two levels of hash tables

Question: Can we make look up time $O(1)$ in worst case?

Perfect Hashing for Static Data

- Do hashing once.
- If $Y_i = |\ell(i)| > 10$ then hash elements of $\ell(i)$ to a table of size Y_i^2 .

Lemma (Look-up)

Expected worst case look up time is $O(1)$.

Lemma (Size)

If $|S| = O(m)$ then space usage of perfect hashing is $O(m)$.

Intuition: Throwing m Balls in to m^2 Bins

- **Pr**[i th ball lands in j th bin]

Intuition: Throwing m Balls in to m^2 Bins

- $\Pr[i\text{th ball lands in } j\text{th bin}] = 1/m^2$
- For a fixed bin j , $Y_j = \#$ balls in bin j .

Intuition: Throwing m Balls in to m^2 Bins

- $\Pr[i\text{th ball lands in } j\text{th bin}] = 1/m^2$
- For a fixed bin j , $Y_j = \#$ balls in bin j . $\mathbf{E}[Y_j] = 1/m$.

Intuition: Throwing m Balls in to m^2 Bins

- $\Pr[i\text{th ball lands in } j\text{th bin}] = 1/m^2$
- For a fixed bin j , $Y_j = \#$ balls in bin j . $\mathbf{E}[Y_j] = 1/m$.
- For $c \geq 3$, let $\delta = cm - 1$. $\Pr[Y_j > c]$?

Intuition: Throwing m Balls in to m^2 Bins

- $\Pr[i\text{th ball lands in } j\text{th bin}] = 1/m^2$
- For a fixed bin j , $Y_j = \#$ balls in bin j . $\mathbf{E}[Y_j] = 1/m$.
- For $c \geq 3$, let $\delta = cm - 1$. $\Pr[Y_j > c]$?

$$\begin{aligned}\Pr[Y_j > cm/m] &= \Pr[Y_j > (1 + \delta) \mathbf{E}[Y_j]] \\ \text{(Chernoff)} &< \left(\frac{e^\delta}{(1+\delta)^{(1+\delta)}} \right)^\mu \\ &= \left(\frac{e^{(cm-1)}}{(cm)^{cm}} \right)^{1/m} \leq (e/c)^c (1/m^c) \\ &\leq 1/m^3\end{aligned}$$

Intuition: Throwing m Balls in to m^2 Bins

- $\Pr[i\text{th ball lands in } j\text{th bin}] = 1/m^2$
- For a fixed bin j , $Y_j = \#$ balls in bin j . $\mathbf{E}[Y_j] = 1/m$.
- For $c \geq 3$, let $\delta = cm - 1$. $\Pr[Y_j > c]$?

$$\begin{aligned}\Pr[Y_j > cm/m] &= \Pr[Y_j > (1 + \delta) \mathbf{E}[Y_j]] \\ \text{(Chernoff)} &< \left(\frac{e^\delta}{(1+\delta)^{(1+\delta)}} \right)^\mu \\ &= \left(\frac{e^{(cm-1)}}{(cm)^{cm}} \right)^{1/m} \leq (e/c)^c (1/m^c) \\ &\leq 1/m^3\end{aligned}$$

- $\Pr\left[\max_{j=1}^{m^2} Y_j > c\right] \leq$

Intuition: Throwing m Balls in to m^2 Bins

- $\Pr[i\text{th ball lands in } j\text{th bin}] = 1/m^2$
- For a fixed bin j , $Y_j = \#$ balls in bin j . $\mathbf{E}[Y_j] = 1/m$.
- For $c \geq 3$, let $\delta = cm - 1$. $\Pr[Y_j > c]$?

$$\begin{aligned}\Pr[Y_j > cm/m] &= \Pr[Y_j > (1 + \delta) \mathbf{E}[Y_j]] \\ (\text{Chernoff}) &< \left(\frac{e^\delta}{(1+\delta)^{(1+\delta)}} \right)^\mu \\ &= \left(\frac{e^{(cm-1)}}{(cm)^{cm}} \right)^{1/m} \leq (e/c)^c (1/m^c) \\ &\leq 1/m^3\end{aligned}$$

- $\Pr\left[\max_{j=1}^{m^2} Y_j > c\right] \leq 1/m$ (Union bound).
- $\Pr\left[\max_{j=1}^{m^2} Y_j \leq c\right] \geq$

Intuition: Throwing m Balls in to m^2 Bins

- $\Pr[i\text{th ball lands in } j\text{th bin}] = 1/m^2$
- For a fixed bin j , $Y_j = \#$ balls in bin j . $\mathbf{E}[Y_j] = 1/m$.
- For $c \geq 3$, let $\delta = cm - 1$. $\Pr[Y_j > c]$?

$$\begin{aligned}\Pr[Y_j > cm/m] &= \Pr[Y_j > (1 + \delta) \mathbf{E}[Y_j]] \\ (\text{Chernoff}) &< \left(\frac{e^\delta}{(1+\delta)^{(1+\delta)}} \right)^\mu \\ &= \left(\frac{e^{(cm-1)}}{(cm)^{cm}} \right)^{1/m} \leq (e/c)^c (1/m^c) \\ &\leq 1/m^3\end{aligned}$$

- $\Pr\left[\max_{j=1}^{m^2} Y_j > c\right] \leq 1/m$ (Union bound).
- $\Pr\left[\max_{j=1}^{m^2} Y_j \leq c\right] \geq 1 - 1/m - (\text{w.h.p.})$
- $\mathbf{E}[\max_j Y_j] \leq$

Intuition: Throwing m Balls in to m^2 Bins

- $\Pr[i\text{th ball lands in } j\text{th bin}] = 1/m^2$
- For a fixed bin j , $Y_j = \#$ balls in bin j . $\mathbf{E}[Y_j] = 1/m$.
- For $c \geq 3$, let $\delta = cm - 1$. $\Pr[Y_j > c]$?

$$\begin{aligned}\Pr[Y_j > cm/m] &= \Pr[Y_j > (1 + \delta) \mathbf{E}[Y_j]] \\ (\text{Chernoff}) &< \left(\frac{e^\delta}{(1+\delta)^{(1+\delta)}} \right)^\mu \\ &= \left(\frac{e^{(cm-1)}}{(cm)^{cm}} \right)^{1/m} \leq (e/c)^c (1/m^c) \\ &\leq 1/m^3\end{aligned}$$

- $\Pr\left[\max_{j=1}^{m^2} Y_j > c\right] \leq 1/m$ (Union bound).
- $\Pr\left[\max_{j=1}^{m^2} Y_j \leq c\right] \geq 1 - 1/m - (\text{w.h.p.})$
- $\mathbf{E}[\max_j Y_j] \leq c + 1 = O(1)$.

Perfect Hashing

Two levels of hash tables

Question: Can we make look up time $O(1)$ in worst case?

Perfect Hashing for Static Data

- Do hashing once.
- If $Y_i = |\ell(i)| > 10$ then hash elements of $\ell(i)$ to a table of size Y_i^2 .

Lemma (Look-up)

Expected worst case look up time is $O(1)$.

Lemma (Size)

If $|S| = O(m)$ then space usage of perfect hashing is $O(m)$.

Perfect Hashing: Proof of Lemma Size

$O(m)$ space usage

h : the primary hash function. $m_i = \# x$ in S such that $h(x) = i$.

Perfect Hashing: Proof of Lemma Size

$O(m)$ space usage

h : the primary hash function. $m_i = \# x$ in S such that $h(x) = i$.

Perfect Hashing: Proof of Lemma Size

$O(m)$ space usage

h : the primary hash function. $m_i = \# x$ in S such that $h(x) = i$.

Claim

$$\mathbf{E} \left[\sum_{i=0}^{m-1} m_i^2 \right] < 3m \text{ where } m = |S|.$$

Perfect Hashing: Proof of Lemma Size

$O(m)$ space usage

h : the primary hash function. $m_i = \# x$ in S such that $h(x) = i$.

Claim

$$\mathbf{E} \left[\sum_{i=0}^{m-1} m_i^2 \right] < 3m \text{ where } m = |S|.$$

Proof.

Let $[h(x) = i]$ represent indicator variable. $m_i = \sum_{x \in S} [h(x) = i]$.

Perfect Hashing: Proof of Lemma Size

$O(m)$ space usage

h : the primary hash function. $m_i = \# x$ in S such that $h(x) = i$.

Claim

$$\mathbf{E} \left[\sum_{i=0}^{m-1} m_i^2 \right] < 3m \text{ where } m = |S|.$$

Proof.

Let $[h(x) = i]$ represent indicator variable. $m_i = \sum_{x \in S} [h(x) = i]$.

$$\sum_i m_i^2 = \sum_i (\sum_{x \in S} [h(x) = i])^2$$

Perfect Hashing: Proof of Lemma Size

$O(m)$ space usage

h : the primary hash function. $m_i = \# x$ in S such that $h(x) = i$.

Claim

$$\mathbf{E} \left[\sum_{i=0}^{m-1} m_i^2 \right] < 3m \text{ where } m = |S|.$$

Proof.

Let $[h(x) = i]$ represent indicator variable. $m_i = \sum_{x \in S} [h(x) = i]$.

$$\begin{aligned} \sum_i m_i^2 &= \sum_i (\sum_{x \in S} [h(x) = i])^2 \\ &= \sum_i (\sum_x [h(x) = i])^2 + 2 \sum_{x < y} [h(x) = i][h(y) = i] \end{aligned}$$

Perfect Hashing: Proof of Lemma Size

$O(m)$ space usage

h : the primary hash function. $m_i = \# x$ in S such that $h(x) = i$.

Claim

$$\mathbf{E} \left[\sum_{i=0}^{m-1} m_i^2 \right] < 3m \text{ where } m = |S|.$$

Proof.

Let $[h(x) = i]$ represent indicator variable. $m_i = \sum_{x \in S} [h(x) = i]$.

$$\begin{aligned} \sum_i m_i^2 &= \sum_i (\sum_{x \in S} [h(x) = i])^2 \\ &= \sum_i (\sum_x [h(x) = i])^2 + 2 \sum_{x < y} [h(x) = i][h(y) = i] \\ &= \sum_x (\sum_i [h(x) = i]) + 2 \sum_{x < y} \sum_i [h(x) = i][h(y) = i] \end{aligned}$$

Perfect Hashing: Proof of Lemma Size

$O(m)$ space usage

h : the primary hash function. $m_i = \# x$ in S such that $h(x) = i$.

Claim

$$\mathbf{E} \left[\sum_{i=0}^{m-1} m_i^2 \right] < 3m \text{ where } m = |S|.$$

Proof.

Let $[h(x) = i]$ represent indicator variable. $m_i = \sum_{x \in S} [h(x) = i]$.

$$\begin{aligned} \sum_i m_i^2 &= \sum_i (\sum_{x \in S} [h(x) = i])^2 \\ &= \sum_i (\sum_x [h(x) = i]^2 + 2 \sum_{x < y} [h(x) = i][h(y) = i]) \\ &= \sum_x (\sum_i [h(x) = i]) + 2 \sum_{x < y} \sum_i [h(x) = i][h(y) = i] \\ &= \sum_x (1) + 2 \sum_{x < y} [h(x) = h(y)] \end{aligned}$$

Perfect Hashing: Proof of Lemma Size

$O(m)$ space usage

h : the primary hash function. $m_i = \# x$ in S such that $h(x) = i$.

Claim

$$\mathbb{E} \left[\sum_{i=0}^{m-1} m_i^2 \right] < 3m \text{ where } m = |S|.$$

Proof.

Let $[h(x) = i]$ represent indicator variable. $m_i = \sum_{x \in S} [h(x) = i]$.

$$\begin{aligned} \sum_i m_i^2 &= \sum_i (\sum_{x \in S} [h(x) = i])^2 \\ &= \sum_i (\sum_x [h(x) = i]^2 + 2 \sum_{x < y} [h(x) = i][h(y) = i]) \\ &= \sum_x (\sum_i [h(x) = i]) + 2 \sum_{x < y} \sum_i [h(x) = i][h(y) = i] \\ &= \sum_x (1) + 2 \sum_{x < y} [h(x) = h(y)] \end{aligned}$$

$$\mathbb{E} [\sum_i m_i^2] = m + 2 \sum_{x < y} \Pr[h(x) = h(y)] = m + 2 \frac{m(m-1)}{2} \frac{1}{m} < 2m$$

Rehashing, amortization and...

... making the hash table dynamic

So far we assumed fixed S of size $\simeq m$.

Question: What happens as items are inserted and deleted?

- 1 If $|S|$ grows to more than cm for some constant c then hash table performance clearly degrades.
- 2 If $|S|$ stays around $\simeq m$ but incurs many insertions and deletions then the initial random hash function is no longer random enough!

Rehashing, amortization and...

... making the hash table dynamic

So far we assumed fixed S of size $\simeq m$.

Question: What happens as items are inserted and deleted?

- 1 If $|S|$ grows to more than cm for some constant c then hash table performance clearly degrades.
- 2 If $|S|$ stays around $\simeq m$ but incurs many insertions and deletions then the initial random hash function is no longer random enough!

Solution: Rebuild hash table periodically!

- 1 Choose a new table size based on current number of elements in the table.
- 2 Choose a new random hash function and rehash the elements.
- 3 Discard old table and hash function.

Question: When to rebuild? How expensive?

Rebuilding the hash table

- 1 Start with table size m where m is some estimate of $|S|$ (can be some large constant).
- 2 If $|S|$ grows to more than twice current table size, build new hash table (choose a new random hash function) with double the current number of elements. Can also use similar trick if table size falls below quarter the size.

Rebuilding the hash table

- 1 Start with table size m where m is some estimate of $|S|$ (can be some large constant).
- 2 If $|S|$ grows to more than twice current table size, build new hash table (choose a new random hash function) with double the current number of elements. Can also use similar trick if table size falls below quarter the size.
- 3 If $|S|$ stays roughly the same but more than $c|S|$ operations on table for some chosen constant c (say **10**), rebuild.

Rebuilding the hash table

- 1 Start with table size m where m is some estimate of $|S|$ (can be some large constant).
- 2 If $|S|$ grows to more than twice current table size, build new hash table (choose a new random hash function) with double the current number of elements. Can also use similar trick if table size falls below quarter the size.
- 3 If $|S|$ stays roughly the same but more than $c|S|$ operations on table for some chosen constant c (say **10**), rebuild.

The **amortize** cost of rebuilding to previously performed operations. Rebuilding ensures $O(1)$ expected analysis holds even when S changes. Hence $O(1)$ expected look up/insert/delete time *dynamic* data dictionary data structure!

Hashing:

- 1 To insert x in dictionary store x in table in location $h(x)$
- 2 To lookup y in dictionary check contents of location $h(y)$
- 3 Storing items in dictionary expensive in terms of memory, especially if items are unwieldy objects such as long strings, images, etc with *non-uniform* sizes.

Bloom Filters

Hashing:

- 1 To insert x in dictionary store x in table in location $h(x)$
- 2 To lookup y in dictionary check contents of location $h(y)$
- 3 Storing items in dictionary expensive in terms of memory, especially if items are unwieldy objects such as long strings, images, etc with *non-uniform* sizes.

Bloom Filter: tradeoff space for false positives

- 1 To insert x in dictionary set *bit* to **1** in location $h(x)$ (initially all bits are set to **0**)
- 2 To lookup y if bit in location $h(y)$ is **1** say yes, else no.

Bloom Filters

Bloom Filter: tradeoff space for false positives

- 1 To insert x in dictionary set *bit* to **1** in location $h(x)$ (initially all bits are set to **0**)
- 2 To lookup y if bit in location $h(y)$ is **1** say yes, else no
- 3 No false negatives but false positives possible due to collisions

Bloom Filters

Bloom Filter: tradeoff space for false positives

- 1 To insert x in dictionary set *bit* to **1** in location $h(x)$ (initially all bits are set to **0**)
- 2 To lookup y if bit in location $h(y)$ is **1** say yes, else no
- 3 No false negatives but false positives possible due to collisions

Reducing false positives:

- 1 Pick k hash functions h_1, h_2, \dots, h_k *independently*

Bloom Filters

Bloom Filter: tradeoff space for false positives

- 1 To insert x in dictionary set *bit* to **1** in location $h(x)$ (initially all bits are set to **0**)
- 2 To lookup y if bit in location $h(y)$ is **1** say yes, else no
- 3 No false negatives but false positives possible due to collisions

Reducing false positives:

- 1 Pick k hash functions h_1, h_2, \dots, h_k *independently*
- 2 To insert set $h_i(x)$ th bit to one in table i for each $1 \leq i \leq k$

Bloom Filters

Bloom Filter: tradeoff space for false positives

- 1 To insert x in dictionary set *bit* to **1** in location $h(x)$ (initially all bits are set to **0**)
- 2 To lookup y if bit in location $h(y)$ is **1** say yes, else no
- 3 No false negatives but false positives possible due to collisions

Reducing false positives:

- 1 Pick k hash functions h_1, h_2, \dots, h_k *independently*
- 2 To insert set $h_i(x)$ th bit to one in table i for each $1 \leq i \leq k$
- 3 To lookup y compute $h_i(y)$ for $1 \leq i \leq k$ and say yes only if each bit in the corresponding location is **1**, otherwise say no. If probability of false positive for one hash function is $\alpha < 1$ then with k independent hash function it is α^k .

Take away points

- ① Hashing is a powerful and important technique for dictionaries. Many practical applications.
- ② Randomization fundamental to understand hashing.
- ③ Good and efficient hashing possible in theory and practice with proper definitions (universal, perfect, etc).
- ④ Related ideas of creating a compact fingerprint/sketch for objects is very powerful in theory and practice.

Practical Issues

Hashing used typically for integers, vectors, strings etc.

- Universal hashing is defined for integers. To implement for other objects need to map objects in some fashion to integers (via representation)
- Practical methods for various important cases such as vectors, strings are studied extensively. See http://en.wikipedia.org/wiki/Universal_hashing for some pointers.
- Details on Cuckoo hashing and its advantage over chaining http://en.wikipedia.org/wiki/Cuckoo_hashing.
- Relatively recent important paper bridging theory and practice of hashing. “The power of simple tabulation hashing” by Mikkel Thorup and Mihai Patrascu, 2011. See http://en.wikipedia.org/wiki/Tabulation_hashing
- Cryptographic hash functions have a different motivation and