

Lecture 5: Hashing

Rafael Oliveira

University of Waterloo
Cheriton School of Computer Science

rafael.oliveira.teaching@gmail.com

May 25, 2021

Overview

- Introduction
 - Hash Functions
 - Why is hashing?
 - How to hash?
- Succinctness of Hash Functions
 - Coping with randomness
 - Universal Hashing
 - Hashing using 2-universal families
 - Perfect Hashing
- Acknowledgements

Computational Model

Before we talk about hash functions, we need to state our model of computation:

Definition (Word RAM model)

In the word RAM^a model:

- all elements are integers that fit in a machine word of w bits
- Basic operations (comparison, arithmetic, bitwise) on such words take $\Theta(1)$ time
- We can also access *any* position in the array in $\Theta(1)$ time

^aRAM stands for Random Access Model

What is hashing?

We want to store n elements (keys) from the set $U = \{0, 1, \dots, m - 1\}$, where $m \gg n$, in a data structure that supports *insertions*, *deletions*, *search* “as efficiently as possible.”

$m < 2^w$ \therefore can store each element of U
in one word

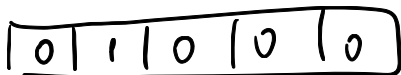
What is hashing?

We want to store n elements (keys) from the set $U = \{0, 1, \dots, m - 1\}$, where $m \gg n$, in a data structure that supports *insertions*, *deletions*, *search* “as efficiently as possible.”

Naive approach: use an array A of m elements, initially $A[i] = 0$ for all i , and when a key is inserted, set $A[i] = 1$.



insert(1)



What is hashing?

We want to store n elements (keys) from the set $U = \{0, 1, \dots, m - 1\}$, where $m \gg n$, in a data structure that supports *insertions*, *deletions*, *search* “as efficiently as possible.”

Naive approach: use an array A of m elements, initially $A[i] = 0$ for all i , and when a key is inserted, set $A[i] = 1$.

- Insertion: $O(1)$, Deletion: $O(1)$, Search: $O(1)$
- Memory: $O(m \log(m))$

(this is very bad!)

if assume $\log(m) < w$
then $O(m)$ space

What is hashing?

We want to store n elements (keys) from the set $U = \{0, 1, \dots, m - 1\}$, where $m \gg n$, in a data structure that supports *insertions*, *deletions*, *search* “as efficiently as possible.”

Naive approach: use an array A of m elements, initially $A[i] = 0$ for all i , and when a key is inserted, set $A[i] = 1$.

- Insertion: $O(1)$, Deletion: $O(1)$, Search: $O(1)$
- Memory: $O(m \log(m))$ (this is very bad!)

Want to also achieve optimal memory $O(n \log(m))$. For this we will use a technique called *hashing*. or be as close as possible

- A *hash function* is a function $h : U \rightarrow [0, n - 1]$, where $|U| = m \gg n$.
- A *hash table* is a data structure that consists of:
 - a table T with n cells $[0, n - 1]$, each cell storing $O(\log(m))$ bits
 - a hash function $h : U \rightarrow [0, n - 1]$

From now on, we will define memory as *# of cells*.

Why is hashing useful?

- Designing efficient data structures (dictionaries) for searching
- Data streaming algorithms
- Derandomization

When randomness needed only involves that pairs, or triples, or *small number of elements* “*look independent.*”

- Cryptography
Construct functions that look random to adversaries, but are easy for us to compute.
- Complexity Theory
- many more

Challenges in Hashing

Setup:

- Universe $U = \{0, \dots, m - 1\}$ of size $m \gg n$ where n is the size of the range of our hash function $h : U \rightarrow [0, n - 1]$
- Store $O(n)$ elements of U (keys) in hash table T (which has n cells)

Challenges in Hashing

Setup:

- Universe $U = \{0, \dots, m - 1\}$ of size $m \gg n$ where n is the size of the range of our hash function $h : U \rightarrow [0, n - 1]$
- Store $O(n)$ elements of U (keys) in hash table T (which has n cells)

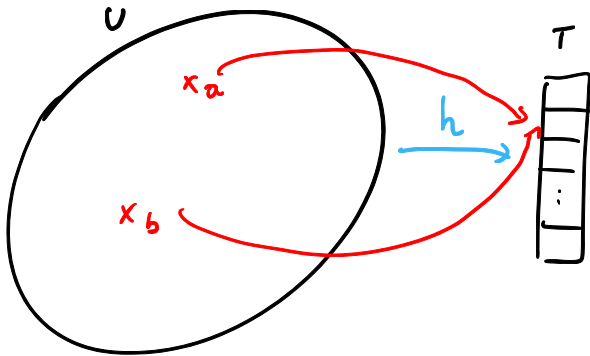
Ideally, want hash function to map *different keys* into *different locations*.

Definition (Collision)

A *collision* happens for hash function h with inputs $x, y \in U$ if

$$x \neq y \text{ and } h(x) = h(y).$$

Pigeonhole principle \Rightarrow impossible without knowing keys in advance.



$$|U| > |T|$$

Challenges in Hashing

Setup:

- Universe $U = \{0, \dots, m - 1\}$ of size $m \gg n$ where n is the size of the range of our hash function $h : U \rightarrow [0, n - 1]$
- Store $O(n)$ elements of U (keys) in hash table T (which has n cells)

Ideally, want hash function to map *different keys* into *different locations*.

Definition (Collision)

A *collision* happens for hash function h with inputs $x, y \in U$ if

$$x \neq y \text{ and } h(x) = h(y).$$

Pigeonhole principle \Rightarrow impossible without knowing keys in advance.

Will settle for: $\#$ collisions *small with high probability*.

Definition (perfect hashing): given a set $S \subset U$
and hash function $h: U \rightarrow [0, \dots, n-1]$

h is perfect for S if $\forall x, y \in S$
 $x \neq y$

$$h(x) \neq h(y)$$

i.e. there are no collisions.

Our solution: family of hash functions

Construct *family* of hash functions \mathcal{H} such that the *number of collisions* is **small** with **high probability**, when we pick hash function uniformly at random from the family \mathcal{H} .

Our solution: family of hash functions

Construct *family* of hash functions \mathcal{H} such that the *number of collisions* is **small** with **high probability**, when we pick hash function uniformly at random from the family \mathcal{H} .

Simplest version to keep in mind:

$$\Pr_{h \in_R \mathcal{H}} [h(x) = h(y)] \leq \frac{1}{\text{poly}(n)} \quad \forall x \neq y \in U$$

Our solution: family of hash functions

Construct *family* of hash functions \mathcal{H} such that the *number of collisions* is **small** with **high probability**, when we pick hash function uniformly at random from the family \mathcal{H} .

Simplest version to keep in mind:

$$\Pr_{h \in \mathcal{H}} [h(x) = h(y)] \leq \frac{1}{\text{poly}(n)} \quad \forall x \neq y \in U$$

Assumptions:

- keys are independent from hash function we choose.
- we **do not** know keys in advance (even if we did, nontrivial problem!)

Question

Still could have collisions. How do we handle them?

Random Hash Functions?

Natural to consider following approach:

From all functions $h : U \rightarrow [0, n - 1]$, just pick one uniformly at random.

Random Hash Functions?

Natural to consider following approach:

From all functions $h : U \rightarrow [0, n - 1]$, just pick one uniformly at random.

This setting is same as our balls-and-bins setting!

$$h : \underbrace{[0, m-1]}_{\text{balls}} \rightarrow \underbrace{[0, n-1]}_{\text{bins}}$$

taking n ^{labelled} balls from U \Leftrightarrow randomly throwing (unlabelled) n balls into bins
(+ h)

Random Hash Functions?

Natural to consider following approach:

From all functions $h : U \rightarrow [0, n - 1]$, just pick one uniformly at random.

This setting is same as our balls-and-bins setting!

So, if we have to store n keys:

- Expected number of keys in a location: 1
- maximum number of collisions (max load) in one particular location: $O(\log n / \log \log n)$ keys

Random Hash Functions?

Natural to consider following approach:

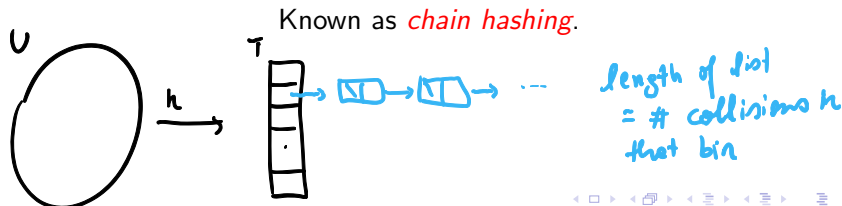
From all functions $h : U \rightarrow [0, n - 1]$, just pick one uniformly at random.

This setting is same as our balls-and-bins setting!

So, if we have to store n keys:

- Expected number of keys in a location: 1
- maximum number of collisions (max load) in one particular location: $O(\log n / \log \log n)$ keys

Solving collisions: store all keys hashed into location i by a linked list.



Random Hash Functions?

Natural to consider following approach:

From all functions $h : U \rightarrow [0, n - 1]$, just pick one uniformly at random.

This setting is same as our balls-and-bins setting!

So, if we have to store n keys:

- Expected number of keys in a location: 1
- maximum number of collisions (max load) in one particular location: $O(\log n / \log \log n)$ keys

Solving collisions: store all keys hashed into location i by a linked list.

Known as *chain hashing*.

Could also pick *two* random hash functions and use *power of two choices*.

Collision bound becomes $O(\log \log n)$.

Random Hash Functions?

Random hash functions look very good. However, we haven't discussed the following:

Random Hash Functions?

Random hash functions look very good. However, we haven't discussed the following:

Question

*How much resource (time & space) does it take to **compute** random hash functions?*

Random Hash Functions?

Random hash functions look very good. However, we haven't discussed the following:

Question

How much resource (time & space) does it take to *compute* random hash functions?

- Storing entire function $h : U \rightarrow [0, n - 1]$ require $O(m \log n)$ bits (way too much space!)
- Even if we only stored the elements we saw, would require $O(n)$ time to evaluate $h(x)$ (need to decide if we had already computed it!)

$$h(1) = 2 \quad h(3) = 1 \quad \boxed{h(17) = 4} \dots$$

17 comes again need to check whether I already defined it
($O(n)$ time)

Random Hash Functions?

Random hash functions look very good. However, we haven't discussed the following:

Question

*How much resource (time & space) does it take to **compute** random hash functions?*

- Storing entire function $h : U \rightarrow [0, n - 1]$ require $O(m \log n)$ bits (way too much space!)
- Even if we only stored the elements we saw, would require $O(n)$ time to evaluate $h(x)$ (need to decide if we had already computed it!)

Remark

Thus, for random function all operations (insert, delete, search) take $O(n \log m)$ time (at best!)

Random Hash Functions?

Random hash functions look very good. However, we haven't discussed the following:

Question

*How much resource (time & space) does it take to **compute** random hash functions?*

- Storing entire function $h : U \rightarrow [0, n - 1]$ require $O(m \log n)$ bits (way too much space!)
- Even if we only stored the elements we saw, would require $O(n)$ time to evaluate $h(x)$ (need to decide if we had already computed it!)

Remark

Thus, for random function all operations (insert, delete, search) take $O(n \log m)$ time (at best!)

How do we cope with the computational problem that arose with randomness?

- Introduction
 - Hash Functions
 - Why is hashing?
 - How to hash?
- Succinctness of Hash Functions
 - Coping with randomness
 - Universal Hashing
 - Hashing using 2-universal families
 - Perfect Hashing
- Acknowledgements

How to cope with “hardness” of randomness?

We want something that is *random-like* (few collisions w.h.p.) but *easy to compute/represent*.

How to cope with “hardness” of randomness?

We want something that is *random-like* (few collisions w.h.p.) but *easy to compute/represent*.

Ideally something that takes $O(\log m)$ time to compute (size of our input).

“word size”
 $O(1)$

How to cope with “hardness” of randomness?

We want something that is *random-like* (few collisions w.h.p.) but *easy to compute/represent*.

Ideally something that takes $O(\log m)$ time to compute (size of our input).

Question

How many hash functions can we have with the property above?

How to cope with “hardness” of randomness?

We want something that is *random-like* (few collisions w.h.p.) but *easy to compute/represent*.

Ideally something that takes $O(\log m)$ time to compute (size of our input).

Question

How many hash functions can we have with the property above?

$\text{poly}(m)$ functions, as each function takes at most $O(\log m)$ bits to describe. Thus these are *succinct functions* (easy to describe and compute) which have *random-like* properties!

$$2^{O(\log m)} = 2^{c \cdot \log m} = m^c = \text{poly}(m)$$

How to cope with “hardness” of randomness?

We want something that is *random-like* (few collisions w.h.p.) but *easy to compute/represent*.

Ideally something that takes $O(\log m)$ time to compute (size of our input).

Question

How many hash functions can we have with the property above?

$\text{poly}(m)$ functions, as each function takes at most $O(\log m)$ bits to describe. Thus these are *succinct functions* (easy to describe and compute) which have *random-like* properties!

Part of *derandomization/pseudorandomness*: huge subfield in TCS!

k-wise independence

Weaker notion of independence.

"sacrifice some randomness to get better computation"

k-wise independence

Weaker notion of independence.

Definition (Full Independence)

A set of random variables X_1, \dots, X_n are said to be (fully) independent if they satisfy

$$\forall a_i \quad \Pr \left[\bigcap_{i=1}^n X_i = a_i \right] = \prod_{i=1}^n \Pr[X_i = a_i]$$

↳ in range of X_i

k-wise independence

Weaker notion of independence.

Definition (Full Independence)

A set of random variables X_1, \dots, X_n are said to be (fully) independent if they satisfy

$$\Pr \left[\bigcap_{i=1}^n X_i = a_i \right] = \prod_{i=1}^n \Pr[X_i = a_i]$$

Definition (k -wise Independence)

A set of random variables X_1, \dots, X_n are said to be k -wise independent if for any set $J \subset [n]$ such that $|J| \leq k$ they satisfy

$\forall a_i$ in range X_i

$$\Pr \left[\bigcap_{i \in J} X_i = a_i \right] = \prod_{i \in J} \Pr[X_i = a_i]$$

small subsets of our random variables

Pairwise independence

When $k = 2$, k -wise independence is called *pairwise independence*.

Example (XOR pairwise independence)

Given b uniformly distributed, independent random bits Y_1, \dots, Y_b , we can generate $2^b - 1$ pairwise independent, uniform random variables as follows:

$$X_S := \bigoplus_{i \in S} Y_i \quad S \subseteq [b] \setminus \emptyset$$

Pairwise independence

When $k = 2$, k -wise independence is called *pairwise independence*.

Example (XOR pairwise independence)

Given b uniformly distributed, independent random bits Y_1, \dots, Y_b , we can generate $2^b - 1$ pairwise independent, uniform random variables as follows:

$$X_S := \bigoplus_{i \in S} Y_i \quad S \subseteq [b] \setminus \emptyset$$

Pairwise independence

When $k = 2$, k -wise independence is called *pairwise independence*.

Example (XOR pairwise independence)

Given b uniformly distributed, independent random bits Y_1, \dots, Y_b , we can generate $2^b - 1$ pairwise independent, uniform random variables as follows:

$$X_S := \bigoplus_{i \in S} Y_i \quad S \subseteq [b] \setminus \emptyset$$

- Why are they uniformly random?

$$S = \{1, 2, 3\}$$

$$X_S = \underbrace{Y_1 \oplus Y_2}_{a} \oplus Y_3$$

\downarrow \downarrow \downarrow
 $\frac{1}{2}$ $\frac{1}{2}$ $\frac{1}{2}$
 a $1-a$

\uparrow
 n
 $\{0, 1\}$

$$\left\{ \begin{array}{l} 0 \text{ w.p. } \frac{1}{2} \\ 1 \text{ w.p. } \frac{1}{2} \end{array} \right.$$

Pairwise independence

When $k = 2$, k -wise independence is called *pairwise independence*.

Example (XOR pairwise independence)

Given b uniformly distributed, independent random bits Y_1, \dots, Y_b , we can generate $2^b - 1$ pairwise independent, uniform random variables as follows:

$$X_S := \bigoplus_{i \in S} Y_i \quad S \subseteq [b] \setminus \emptyset$$

- Why are they uniformly random?
- Why are they pairwise independent?

$$\begin{aligned} & \Pr[X_{S_1} = a_1 \text{ and } X_{S_2} = a_2] \\ &= \Pr[X_{S_1} = a_1] \cdot \Pr[X_{S_2} = a_2] \end{aligned}$$

$$S_1 \neq S_2 \quad \exists j \in S_2 \setminus S_1$$

even if we set all variables in S_1

$$X_{S_2} = Y_j \oplus \boxed{}$$

$\xrightarrow{\text{not set}}$ $\xrightarrow{k \neq 0}$ $\xrightarrow{1/2}$

Pairwise independence

When $k = 2$, k -wise independence is called *pairwise independence*.

Example (XOR pairwise independence)

Given b uniformly distributed, independent random bits Y_1, \dots, Y_b , we can generate $2^b - 1$ pairwise independent, uniform random variables as follows:

$$X_S := \bigoplus_{i \in S} Y_i \quad S \subseteq [b] \setminus \emptyset$$

- Why are they uniformly random?
- Why are they pairwise independent?
- Are they also 3-wise independent?

$$X_{\{1\}} \quad X_{\{2\}} \quad X_{\{1,2\}}$$

0 0 \Rightarrow 0 u.p. 1

Pairwise independence II

Example (Pairwise independence in \mathbb{F}_p)

Let p be a prime number. Given 2 uniformly random variables $Y_1, Y_2 \sim [0, \dots, p-1]$, generate p pairwise independent random variables as follows:

$$X_i := Y_1 + i \cdot Y_2 \pmod{p} \quad i \in [0, p-1]$$

$$\mathbb{F}_p = \mathbb{Z}/p\mathbb{Z} \quad (\text{integers mod } p)$$

Pairwise independence II

Example (Pairwise independence in \mathbb{F}_p)

Let p be a prime number. Given 2 uniformly random variables $Y_1, Y_2 \sim [0, \dots, p-1]$, generate p pairwise independent random variables as follows:

$$X_i := Y_1 + i \cdot Y_2 \pmod{p} \quad i \in [0, p-1]$$

- Why are they uniformly random?

$$X_i = Y_1 + Y_2 \pmod{p}$$

$$a \in \mathbb{F}_p \quad (x, a-x)$$

have p pairs (Y_1, Y_2)

$$\text{s.t. } X_i(Y_1, Y_2) = a$$

$$\Pr[X_i = a] = \frac{p}{p^2} = \frac{1}{p}$$

Pairwise independence II

Example (Pairwise independence in \mathbb{F}_p)

Let p be a prime number. Given 2 uniformly random variables $Y_1, Y_2 \sim [0, \dots, p-1]$, generate p pairwise independent random variables as follows:

$$X_i := Y_1 + i \cdot Y_2 \pmod{p} \quad i \in [0, p-1]$$

- Why are they uniformly random?
- Why are they pairwise independent?

$$P_{\mathcal{X}}[X_i = a_i \wedge X_j = a_j] = \frac{1}{p^2}$$

exactly independent

$$\left. \begin{array}{l} Y_1 + i Y_2 = a_i \pmod{p} \\ Y_1 + j Y_2 = a_j \pmod{p} \end{array} \right\} \begin{array}{l} \text{system of 2 equations} \\ \text{in 2 variables} \end{array} \quad (\text{independent})$$

\therefore we have exactly one solution

$$(j-i)Y_2 = (a_j - a_i) \pmod{p}$$

$\neq 0$

$$\therefore P_{\mathcal{X}}[X_i = a_i \wedge X_j = a_j] = \frac{1}{p^2}$$

Pairwise independence II

Example (Pairwise independence in \mathbb{F}_p)

Let p be a prime number. Given 2 uniformly random variables $Y_1, Y_2 \sim [0, \dots, p-1]$, generate p pairwise independent random variables as follows:

$$X_i := Y_1 + i \cdot Y_2 \pmod{p} \quad i \in [0, p-1]$$

- Why are they uniformly random?
- Why are they pairwise independent?
- Are they also 3-wise independent?

NO.

$$X_0 = 0$$

$$X_1 = 0$$

$$X_2 = 1$$

Prob. 0 of happening

(inconsistent system of equations)

Pairwise independence II

Example (Pairwise independence in \mathbb{F}_p)

Let p be a prime number. Given 2 uniformly random variables $Y_1, Y_2 \sim [0, \dots, p-1]$, generate p pairwise independent random variables as follows:

$$X_i := Y_1 + i \cdot Y_2 \pmod{p} \quad i \in [0, p-1]$$

- Why are they uniformly random?
- Why are they pairwise independent?
- Are they also 3-wise independent?

Can think of these random variables as picking a random line.

- If we only know one point of the line, the second point is still uniformly random. *pairwise independence*
- Two points determine the line. *not 3-wise independent*

- Introduction
 - Hash Functions
 - Why is hashing?
 - How to hash?
- Succinctness of Hash Functions
 - Coping with randomness
 - **Universal Hashing**
 - Hashing using 2-universal families
 - Perfect Hashing
- Acknowledgements

Universal Hash Functions

We want hash functions. Why are we talking about random variables?

Universal Hash Functions

We want hash functions. Why are we talking about random variables?

Definition (Universal Hash Functions)

Let U be a universe with $|U| \geq n$. A family of hash functions $\mathcal{H} = \{h : U \rightarrow [0, n - 1]\}$ is *k-universal* if, *for any distinct* elements $u_1, \dots, u_k \in U$, we have

$$\Pr_{h \in \mathcal{H}} [h(u_1) = h(u_2) = \dots = h(u_k)] \leq 1/n^{k-1}$$

Universal Hash Functions

We want hash functions. Why are we talking about random variables?

Definition (Universal Hash Functions)

Let U be a universe with $|U| \geq n$. A family of hash functions $\mathcal{H} = \{h : U \rightarrow [0, n - 1]\}$ is *k-universal* if, *for any distinct* elements $u_1, \dots, u_k \in U$, we have

$$\Pr_{h \in_R \mathcal{H}} [h(u_1) = h(u_2) = \dots = h(u_k)] \leq 1/n^{k-1}$$

Definition (Strongly Universal Hash Functions)

$\mathcal{H} = \{h : U \rightarrow [0, n - 1]\}$ is *strongly k-universal* if, *for any distinct* elements $u_1, \dots, u_k \in U$ and *for any* values $y_1, \dots, y_k \in [0, n - 1]$, we have

$$\Pr_{h \in_R \mathcal{H}} [h(u_1) = y_1, \dots, h(u_k) = y_k] \leq 1/n^k$$

Relation to k -wise independent random variables

What do the previous definitions have to do with random variables?

Relation to k -wise independent random variables

What do the previous definitions have to do with random variables?

Family \mathcal{H} is *strongly k -universal* if the random variables $h(0), \dots, h(|U| - 1)$ are *k -wise independent*.

Relation to k -wise independent random variables

What do the previous definitions have to do with random variables?

Family \mathcal{H} is *strongly k -universal* if the random variables $h(0), \dots, h(|U| - 1)$ are *k -wise independent*.

Can use random variables to construct universal hash functions!

Strongly 2-universal families of hash functions

Let p be a prime number, $U = [0, p - 1]$.

Proposition

$$\mathcal{H} = \{h_{a,b}(x) := a \cdot x + b \pmod{p} \mid a, b \in [0, p - 1]\}$$

is strongly 2-universal.

$$\begin{array}{cc} h_{a,b}(i) & h_{a,b}(j) \\ \parallel & \parallel \\ a \cdot i + b & a \cdot j + b \end{array}$$

Same proof we did before.

Strongly 2-universal families of hash functions

Let p be a prime number, $U = [0, p - 1]$.

Proposition

$$\mathcal{H} = \{h_{a,b}(x) := a \cdot x + b \pmod p \mid a, b \in [0, p - 1]\}$$

is strongly 2-universal.

How do we make the domain U much larger than image of the maps? (as usually in hashing size of universe much larger than size of table)

Strongly 2-universal families of hash functions

Let p be a prime number, $U = [0, p - 1]$.

Proposition

$$\mathcal{H} = \{h_{a,b}(x) := a \cdot x + b \pmod{p} \mid a, b \in [0, p - 1]\}$$

is strongly 2-universal.

How do we make the domain U much larger than image of the maps? (as usually in hashing size of universe much larger than size of table)

Proposition

Let $U = [0, p^k - 1] \equiv [0, p - 1]^k$ and $a = (a_0, \dots, a_{k-1}) \in \mathbb{F}_p^k$

$$\mathcal{H} = \{h_{a,b}(x) := \langle a, x \rangle + b \pmod{p} \mid a \in U, b \in [0, p - 1]\}$$

is strongly 2-universal.

$$\hookrightarrow a_0 x_0 + a_1 x_1 + \dots + a_{k-1} x_{k-1}$$

Strongly 2-universal families of hash functions

Proposition

Let $U = [0, p^k - 1] \equiv [0, p - 1]^k$ and $a = (a_0, \dots, a_{k-1})$

$$\mathcal{H} = \{h_{a,b}(x) := \langle a, x \rangle + b \pmod{p} \mid a \in U, b \in [0, p - 1]\}$$

is strongly 2-universal.

Proof similar to first proof (subsets)
(there $p=2$)

2-universal families of hash functions

$$\mathbb{F}_p \rightarrow \mathbb{F}_p \rightarrow \mathbb{Z}/n\mathbb{Z} \quad \text{simple hash function}$$

\cup

What if my hash table size is not a prime?

Proposition

$$\mathcal{H} = \{h_{a,b}(x) := ((a \cdot x + b) \bmod p) \bmod n \mid a, b \in [0, p-1], a \neq 0\}$$

is 2-universal (but not strongly 2-universal).

Practice problem: prove the proposition above.

k -universal families of hash functions

Can we construct k -universal families of hash functions like this?

k -universal families of hash functions

Can we construct k -universal families of hash functions like this?

- YES! Instead of constructing random lines (degree 1 polynomials), can construct random univariate polynomials of degree $k - 1$

$$h_{a,b}(x) = ax + b$$

$$h_{a,b,c}(x) = ax^2 + bx + c \quad \text{3-wise independent}$$

k -universal families of hash functions

Can we construct k -universal families of hash functions like this?

- YES! Instead of constructing random lines (degree 1 polynomials), can construct random univariate polynomials of degree $k - 1$
- Two points determine a line. Similarly, k points determine a univariate polynomial of degree $k - 1$

k -universal families of hash functions

Can we construct k -universal families of hash functions like this?

- YES! Instead of constructing random lines (degree 1 polynomials), can construct random univariate polynomials of degree $k - 1$
- Two points determine a line. Similarly, k points determine a univariate polynomial of degree $k - 1$
- Random degree $k - 1$ polynomials are k -wise independent!
- Practice problem: prove this!

Efficiency

How did pairwise independence improve over random functions?

Efficiency

How did pairwise independence improve over random functions?

Remark

For random function all operations (insert, delete, search) take $O(n \log m)$ time (at best!) to compute. And it takes $O(n \log m)$ storage!

Efficiency

How did pairwise independence improve over random functions?

Remark

For random function all operations (insert, delete, search) take $O(n \log m)$ time (at best!) to compute. And it takes $O(n \log m)$ storage!

Remark

- In XOR example, our function takes $O(b)$ storage space, and $O(b)$ time to compute.^a *store $x+1$ S get b bits evaluate $\oplus Y_i$ ies $O(b \log b)$*
- In \mathbb{F}_p examples, our function takes $O(1)$ storage space and $O(1)$ time to compute!^b *store a, b defining $h_{a,b}$ compute $x \mapsto ax+b$*

^aReminder that we assume that $b < w$.

^bWe assume that $p < 2^w$.

- Introduction
 - Hash Functions
 - Why is hashing?
 - How to hash?
- Succinctness of Hash Functions
 - Coping with randomness
 - Universal Hashing
 - Hashing using 2-universal families
 - Perfect Hashing
- Acknowledgements

Hashing with 2-universal families

- Let $U = [0, m - 1]$, and p be a prime number such that $m \leq p < 2m$ (exists by Bertrand's postulate)

Hashing with 2-universal families

- Let $U = [0, m - 1]$, and p be a prime number such that $m \leq p < 2m$ (exists by Bertrand's postulate)
- $\mathcal{H} = \{h_{a,b}(x) := ((a \cdot x + b) \bmod p) \bmod n \mid a, b \in [0, p - 1], a \neq 0\}$

Hashing with 2-universal families

- Let $U = [0, m - 1]$, and p be a prime number such that $m \leq p < 2m$ (exists by Bertrand's postulate)
- $\mathcal{H} = \{h_{a,b}(x) := ((a \cdot x + b) \bmod p) \bmod m \mid a, b \in [0, p - 1], a \neq 0\}$
- Only need $0 \neq a$ and $b \in [0, p - 1]$ to store a function from \mathcal{H} .

$O(\log m)$ storage space

Hashing with 2-universal families

- Let $U = [0, m - 1]$, and p be a prime number such that $m \leq p < 2m$ (exists by Bertrand's postulate)
- $\mathcal{H} = \{h_{a,b}(x) := ((a \cdot x + b) \bmod p) \bmod m \mid a, b \in [0, p - 1], a \neq 0\}$
- Only need $0 \neq a$ and $b \in [0, p - 1]$ to store a function from \mathcal{H} .
- Computation time of $h_{a,b}$ is also $O(\log m)$

Hashing with 2-universal families

- Let $U = [0, m - 1]$, and p be a prime number such that $m \leq p < 2m$ (exists by Bertrand's postulate)
- $\mathcal{H} = \{h_{a,b}(x) := ((a \cdot x + b) \bmod p) \bmod n \mid a, b \in [0, p - 1], a \neq 0\}$
- Only need $0 \neq a$ and $b \in [0, p - 1]$ to store a function from \mathcal{H} .
- Computation time of $h_{a,b}$ is also $O(\log m)$
- Can this hash function match chain hashing parameters?
($O(\log \log n)$ search time)

No.

Hashing with 2-universal families

- Let $U = [0, m - 1]$, and p be a prime number such that $m \leq p < 2m$ (exists by Bertrand's postulate)
- $\mathcal{H} = \{h_{a,b}(x) := ((a \cdot x + b) \bmod p) \bmod m \mid a, b \in [0, p - 1], a \neq 0\}$
- Only need $0 \neq a$ and $b \in [0, p - 1]$ to store a function from \mathcal{H} .
- Computation time of $h_{a,b}$ is also $O(\log m)$
- Can this hash function match chain hashing parameters?
($O(\log \log n)$ search time)

Do not have same expected search time as chain hashing.

Hashing with 2-universal families

- Let $U = [0, m - 1]$, and p be a prime number such that $m \leq p < 2m$ (exists by Bertrand's postulate)
- $\mathcal{H} = \{h_{a,b}(x) := ((a \cdot x + b) \bmod p) \bmod n \mid a, b \in [0, p - 1], a \neq 0\}$
- Only need $0 \neq a$ and $b \in [0, p - 1]$ to store a function from \mathcal{H} .
- Computation time of $h_{a,b}$ is also $O(\log m)$
- Can this hash function match chain hashing parameters?
($O(\log \log n)$ search time)

↪ using two random hash functions
Do not have same expected search time as chain hashing.

Lemma (Maximum number of collisions)

The expected number of collisions when hashing ℓ elements using a 2-universal hash family is

$$\leq \ell^2 / 2n$$

Hashing with 2-universal families

Lemma (Maximum number of collisions)

The expected number of collisions ^{pairs} when hashing l elements using a 2-universal hash family is

$$\approx l^2/2n$$

$X_{ij} = \begin{cases} 1 & \text{if keys } i \text{ and } j \text{ are mapped to same location} \\ 0 & \text{otherwise} \end{cases}$

collision btw i,j

$$X = \sum_{1 \leq i < j \leq l} X_{ij}$$

total # collisions

$$E[X] = \sum_{1 \leq i < j \leq l} E[X_{ij}] = \sum_{1 \leq i < j \leq l} P_k [h(i) = h(j)]$$

$\leq \frac{1}{n}$
 \downarrow
2-universal

$$\leq \binom{l}{2} \cdot \frac{1}{n} \approx \frac{l^2}{2n}$$

$\boxed{1, 2, 3, 4, \dots, l}$

$\rightarrow k^2$ collision pairs \checkmark load is k
load is square root of # collisions

Hashing with 2-universal families

Lemma (Maximum number of collisions)

The expected number of collisions when hashing ℓ elements using a 2-universal hash family is

$$\ell^2/2n$$

Thus, by Markov's inequality, we have $P_x[X \geq t] \leq \frac{E[X]}{t}$

Lemma (Maximum load of entry of hash table)

With probability $\geq 1/2$ the ~~number of collisions~~ **max load** when hashing ℓ elements using a 2-universal hash family is

$$\leq \sqrt{\frac{2\ell^2}{n}}$$

When $\ell \approx n$ (as is usually assumed in hashing), we expect $\sqrt{2n}$.

Can prove max load is $\sqrt{\max \# \text{ collisions}}$

- Introduction
 - Hash Functions
 - Why is hashing?
 - How to hash?
- Succinctness of Hash Functions
 - Coping with randomness
 - Universal Hashing
 - Hashing using 2-universal families
 - Perfect Hashing
- Acknowledgements

Perfect Hashing

- **Setup:** we are given *in advance* a static set $S \subset U$ of size n
- How to build a hash table with $O(1)$ search time and $O(n)$ memory?
- Can we still do it with a 2-universal family of hash functions?

Perfect Hashing

- **Setup:** we are given *in advance* a static set $S \subset U$ of size n
- How to build a hash table with $O(1)$ search time and $O(n)$ memory?
- Can we still do it with a 2-universal family of hash functions?

Corollary

If $h \in \mathcal{H}$ is a random hash function from a 2-universal family of hash functions, then for any set $S \subseteq U$ of size $\ell \leq \sqrt{n}$, the probability of h being perfect for S is at least $1/2$.

Proof: There is no collision with probability $\geq 1/2$.

Perfect Hashing

- **Setup:** we are given *in advance* a static set $S \subset U$ of size n
- How to build a hash table with $O(1)$ search time and $O(n)$ memory?
- Can we still do it with a 2-universal family of hash functions?

Corollary

If $h \in \mathcal{H}$ is a random hash function from a 2-universal family of hash functions, then for any set $S \subseteq U$ of size $\ell \leq \sqrt{n}$, the probability of h being perfect for S is at least $1/2$.

Proof: There is no collision with probability $\geq 1/2$.

New idea: build a *two-level* hash table!

Perfect Hashing

- **Setup:** we are given *in advance* a static set $S \subset U$ of size n
- How to build a hash table with $O(1)$ search time and $O(n)$ memory?
- Can we still do it with a 2-universal family of hash functions?

Corollary

If $h \in \mathcal{H}$ is a random hash function from a 2-universal family of hash functions, then for any set $S \subseteq U$ of size $\ell \leq \sqrt{n}$, the probability of h being perfect for S is at least $1/2$.

Proof: There is no collision with probability $\geq 1/2$.

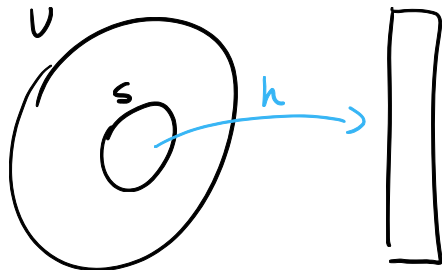
New idea: build a *two-level* hash table!

Theorem

The two-level approach gives perfect hashing scheme.

Proof (sketch) of Theorem

- Pick first layer hash function h uniformly at random from our 2-universal family \mathcal{H} . Test h on our set S .

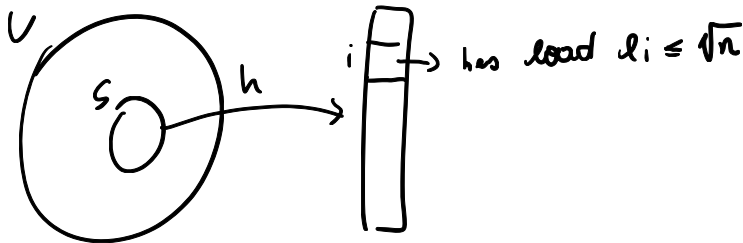


Proof (sketch) of Theorem

- Pick first layer hash function h uniformly at random from our 2-universal family \mathcal{H} . Test h on our set S .
- With probability $\geq 1/2$, max number of collisions in one bin is $\leq \sqrt{n}$. Thus, we will find good hash function for first layer with constant many tries. (with high probability)

Proof (sketch) of Theorem

- Pick first layer hash function h uniformly at random from our 2-universal family \mathcal{H} . Test h on our set S .
- With probability $\geq 1/2$, max number of collisions in one bin is $\leq \sqrt{n}$. Thus, we will find good hash function for first layer with constant many tries. (with high probability)
- Assume max number of collisions (h, S) is $\leq \sqrt{n}$. Let l_i be the load at i^{th} cell of hash table given by h



Proof (sketch) of Theorem

- Pick first layer hash function h uniformly at random from our 2-universal family \mathcal{H} . Test h on our set S .
- With probability $\geq 1/2$, max number of collisions in one bin is $\leq \sqrt{n}$. Thus, we will find good hash function for first layer with constant many tries. (with high probability)
- Assume max number of collisions (h, S) is $\leq \sqrt{n}$. Let l_i be the load at i^{th} cell of hash table given by h
- Thus $l_i \leq \sqrt{n}$ and $\sum_{i=1}^n l_i = |S| = n$

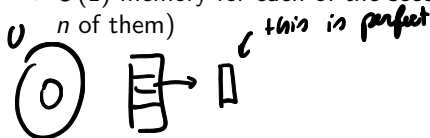
Proof (sketch) of Theorem

- Pick first layer hash function h uniformly at random from our 2-universal family \mathcal{H} . Test h on our set S .
- With probability $\geq 1/2$, max number of collisions in one bin is $\leq \sqrt{n}$. Thus, we will find good hash function for first layer with constant many tries. (with high probability)
- Assume max number of collisions (h, S) is $\leq \sqrt{n}$. Let l_i be the load at i^{th} cell of hash table given by h
- Thus $l_i \leq \sqrt{n}$ and $\sum_{i=1}^n l_i = |S| = n$
- By our lemma, if take $h_i : S \rightarrow [l_i^2]$ from our 2-universal hash family, h_i is perfect with high probability

for S_i
↑
set elements mapped to i
($|S_i| = l_i \leq \sqrt{n}$)

Proof (sketch) of Theorem

- Pick first layer hash function h uniformly at random from our 2-universal family \mathcal{H} . Test h on our set S .
- With probability $\geq 1/2$, max number of collisions in one bin is $\leq \sqrt{n}$. Thus, we will find good hash function for first layer with constant many tries. (with high probability)
- Assume max number of collisions (h, S) is $\leq \sqrt{n}$. Let l_i be the load at i^{th} cell of hash table given by h
- Thus $l_i \leq \sqrt{n}$ and $\sum_{i=1}^n l_i = |S| = n$
- By our lemma, if take $h_i : S \rightarrow [l_i^2]$ from our 2-universal hash family, h_i is perfect with high probability
- **Memory needed:** to store all hash functions, need $O(n)$ memory
 - $O(1)$ memory for the top level hash function h
 - $O(1)$ memory for each of the second-layer hash functions (and we have



Proof (sketch) of Theorem

- Pick first layer hash function h uniformly at random from our 2-universal family \mathcal{H} . Test h on our set S .
- With probability $\geq 1/2$, max number of collisions in one bin is $\leq \sqrt{n}$. Thus, we will find good hash function for first layer with constant many tries. (with high probability)
- Assume max number of collisions (h, S) is $\leq \sqrt{n}$. Let ℓ_i be the load at i^{th} cell of hash table given by h
- Thus $\ell_i \leq \sqrt{n}$ and $\sum_{i=1}^n \ell_i = |S| = n$
- By our lemma, if take $h_i : S \rightarrow [\ell_i^2]$ from our 2-universal hash family, h_i is perfect with high probability
- **Memory needed:** to store all hash functions, need $O(n)$ memory
 - $O(1)$ memory for the top level hash function h
 - $O(1)$ memory for each of the second-layer hash functions (and we have n of them)
- **Time to hash:** $O(1)$ time to evaluate each hash function, and we only have two layers. So total time $O(1)$


Acknowledgement

- Lecture based largely on Lap Chi's notes and on [CLRS 2009, Chapter 11].
- See Lap Chi's notes at <https://cs.uwaterloo.ca/~lapchi/cs466/notes/L05.pdf>

References I

 Motwani, Rajeev and Raghavan, Prabhakar (2007)
Randomized Algorithms

 Mitzenmacher, Michael, and Eli Upfal (2017)
Probability and computing: Randomization and probabilistic techniques in algorithms and data analysis.
Cambridge university press, 2017.

 Cormen, Thomas and Leiserson, Charles and Rivest, Ronald and Stein, Clifford.
(2009)
Introduction to Algorithms, third edition.
MIT Press