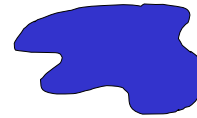


Hashing

CSE 373
Data Structures and Algorithms

Hash Tables

- Constant time accesses!
- A **hash table** is an array of some fixed size, usually a prime number.
- General idea:



key space (e.g., integers, strings)

hash function:
 $h(K)$



hash table



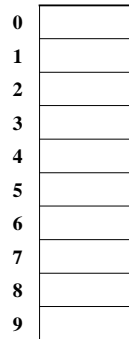
04/15/2009

Hashing

2

Example

- key space = integers
- TableSize = 10
- $h(K) = K \bmod 10$
- **Insert:** 7, 18, 41, 94



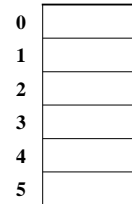
04/15/2009

Hashing

3

Another Example

- key space = integers
- TableSize = 6
- $h(K) = K \bmod 6$
- **Insert:** 7, 18, 41, 34



Student Activity

Hashing

4

Hash Functions

1. **simple/fast** to compute,
2. Avoid **collisions**
3. have keys distributed **evenly** among cells.

Perfect Hash function:

04/15/2009

Hashing

5

Sample Hash Functions:

- key space = strings
 - $s = s_0 s_1 s_2 \dots s_{k-1}$
1. $h(s) = s_0 \bmod \text{TableSize}$
 2. $h(s) = \left(\sum_{i=0}^{k-1} s_i \right) \bmod \text{TableSize}$
 3. $h(s) = \left(\sum_{i=0}^{k-1} s_i \cdot 37^i \right) \bmod \text{TableSize}$

04/15/2009

Hashing

6

Designing a Hash Function for web URLs

$$s = s_0 s_1 s_2 \dots s_{k-1}$$

Issues to take into account:

$$h(s) =$$

Student Activity

Hashing

7

Compare Hash Functions for web URLs

$$h(s) = \left(\sum_{i=0}^n s.charAt(i) \right) \text{ mod TableSize} \quad h(s) = \left(\sum_{i=0}^n s.charAt(i) \cdot 39^i \right) \text{ mod TableSize}$$

Pros:

Pros:

Cons:

Cons:

04/15/2009

Hashing

8

Collision Resolution

Collision: when two keys map to the same location in the hash table.

Two ways to resolve collisions:

1. Separate Chaining
2. Open Addressing (linear probing, quadratic probing, double hashing)

04/15/2009

Hashing

9

Separate Chaining

0	
1	
2	
3	
4	
5	
6	
7	
8	
9	

Insert:
10
22
107
12
42

- **Separate chaining:**
All keys that map to the same hash value are kept in a list ("bucket").

04/15/2009

Hashing

10

Analysis of find

- **Defn:** The **load factor**, λ , of a hash table is the ratio: $\frac{N}{M}$ ← no. of elements
 M ← table size

For separate chaining, λ = average # of elements in a bucket

- unsuccessful:
- successful:

04/15/2009

Hashing

11

How big should the hash table be?

- For Separate Chaining:

04/15/2009

Hashing

12

tableSize: Why Prime?

- Suppose
 - data stored in hash table: 7160, 493, 60, 55, 321, 900, 810
 - tableSize = 10
data hashes to 0, 3, 0, 5, 1, 0, 0
 - tableSize = 11
data hashes to 10, 9, 5, 0, 2, 9, 7

Real-life data tends to have a pattern

Being a multiple of 11 is usually *not* the pattern ☹

04/15/2009 Hashing 13

Open Addressing

0		Insert: 38 19 8 109 10
1		
2		
3		
4		
5		
6		
7		
8		
9		

- **Linear Probing:** after checking spot $h(k)$, try spot $h(k)+1$, if that is full, try $h(k)+2$, then $h(k)+3$, etc.

04/15/2009 Hashing 14

Terminology Alert!

“Open Hashing”	“Closed Hashing”
equals	equals
“Separate Chaining”	“Open Addressing”

Weiss

04/15/2009 Hashing 15

Linear Probing

- $f(i) = i$
- Probe sequence:
 - 0th probe = $h(k) \bmod \text{TableSize}$
 - 1th probe = $(h(k) + 1) \bmod \text{TableSize}$
 - 2th probe = $(h(k) + 2) \bmod \text{TableSize}$
 - ...
 - ith probe = $(h(k) + i) \bmod \text{TableSize}$

04/15/2009 Hashing 16

Write pseudocode for find(k) for Open Addressing with linear probing

- Find(k) returns i where $T(i) = k$

Student Activity

04/15/2009 Hashing 17

Linear Probing – Clustering

[R. Sedgwick]

04/15/2009 Hashing 18

Load Factor in Linear Probing

- For any $\lambda < 1$, linear probing will find an empty slot
- Expected # of probes (for large table sizes)
 - successful search: $\frac{1}{2} \left(1 + \frac{1}{1-\lambda} \right)$
 - unsuccessful search: $\frac{1}{2} \left(1 + \frac{1}{(1-\lambda)^2} \right)$
- Linear probing suffers from *primary clustering*
- Performance quickly degrades for $\lambda > 1/2$

04/15/2009

Hashing

19

Quadratic Probing

$$f(i) = i^2$$

Less likely to encounter Primary Clustering

- Probe sequence:
 - 0th probe = $h(k) \bmod \text{TableSize}$
 - 1th probe = $(h(k) + 1) \bmod \text{TableSize}$
 - 2th probe = $(h(k) + 4) \bmod \text{TableSize}$
 - 3th probe = $(h(k) + 9) \bmod \text{TableSize}$
 - ...
 - i^{th} probe = $(h(k) + i^2) \bmod \text{TableSize}$

04/15/2009

Hashing

20

Quadratic Probing

0	
1	
2	
3	
4	
5	
6	
7	
8	
9	

Insert:
89
18
49
58
79

04/15/2009

Hashing

21

Quadratic Probing Example

insert(76) insert(40) insert(48) insert(5) insert(55)
 $76\%7 = 6$ $40\%7 = 5$ $48\%7 = 6$ $5\%7 = 5$ $55\%7 = 6$

0	
1	
2	
3	
4	
5	
6	76

But... insert(47)
 $47\%7 = 5$

04/15/2009

Hashing

22

Quadratic Probing: Success guarantee for $\lambda < 1/2$

- If size is prime and $\lambda < 1/2$, then quadratic probing will find an empty slot in size/2 probes or fewer.
 - show for all $0 \leq i, j \leq \text{size}/2$ and $i \neq j$

$$(h(x) + i^2) \bmod \text{size} \neq (h(x) + j^2) \bmod \text{size}$$
 - by contradiction: suppose that for some $i \neq j$:

$$(h(x) + i^2) \bmod \text{size} = (h(x) + j^2) \bmod \text{size}$$

$$\Rightarrow i^2 \bmod \text{size} = j^2 \bmod \text{size}$$

$$\Rightarrow (i^2 - j^2) \bmod \text{size} = 0$$

$$\Rightarrow [(i + j)(i - j)] \bmod \text{size} = 0$$
 BUT size does not divide $(i-j)$ or $(i+j)$

04/15/2009

Hashing

23

Quadratic Probing: Properties

- For any $\lambda < 1/2$, quadratic probing will find an empty slot; for bigger λ , quadratic probing may find a slot
- Quadratic probing does not suffer from *primary clustering*: keys hashing to the same *area* are not bad
- But what about keys that hash to the same *spot*?
 - *Secondary Clustering!*

04/15/2009

Hashing

24

Quadratic Probing Works for $\lambda < 1/2$

- If HSize is prime then $(h(x) + i^2) \bmod \text{HSize} \neq (h(x) + j^2) \bmod \text{HSize}$ for $i \neq j$ and $0 \leq i, j < \text{HSize}/2$.
- Proof
 - $(h(x) + i^2) \bmod \text{HSize} = (h(x) + j^2) \bmod \text{HSize}$
 - $(h(x) + i^2) - (h(x) + j^2) \bmod \text{HSize} = 0$
 - $(i^2 - j^2) \bmod \text{HSize} = 0$
 - $(i-j)(i+j) \bmod \text{HSize} = 0$
 - $\Rightarrow \Leftarrow$ HSize does not divide $(i-j)$ or $(i+j)$

04/15/2009

Hashing

25

Double Hashing

$$f(i) = i * g(k)$$

where g is a second hash function

- Probe sequence:
 - 0^{th} probe = $h(k) \bmod \text{TableSize}$
 - 1^{th} probe = $(h(k) + g(k)) \bmod \text{TableSize}$
 - 2^{th} probe = $(h(k) + 2 * g(k)) \bmod \text{TableSize}$
 - 3^{th} probe = $(h(k) + 3 * g(k)) \bmod \text{TableSize}$
 - ...
 - i^{th} probe = $(h(k) + i * g(k)) \bmod \text{TableSize}$

04/15/2009

Hashing

26

Double Hashing Example

$$h(k) = k \bmod 7 \text{ and } g(k) = 5 - (k \bmod 5)$$

	76	93	40	47	10	55
0						
1				47	47	47
2		93	93	93	93	93
3					10	10
4						55
5			40	40	40	40
6	76	76	76	76	76	76
Probes	1	1	1	2	1	2

04/15/2009

Hashing

27

Resolving Collisions with Double Hashing

0	
1	
2	
3	
4	
5	
6	
7	
8	
9	

Hash Functions:
 $H(k) = k \bmod M$
 $H_2(k) = 1 + ((k/M) \bmod (M-1))$
 $M =$

Insert these values into the hash table in this order. Resolve any collisions with double hashing:
 13
 28
 33
 147
 43

04/15/2009

Hashing

28

Rehashing

Idea: When the table gets too full, create a bigger table (usually 2x as large) and hash all the items from the original table into the new table.

- When to rehash?
 - half full ($\lambda = 0.5$)
 - when an insertion fails
 - some other threshold
- Cost of rehashing?

04/15/2009

Hashing

29

Hashing Summary

- Hashing is one of the most important data structures.
- Hashing has many applications where operations are limited to find, insert, and delete.
- Dynamic hash tables have good amortized complexity.

04/15/2009

Hashing

30