

Hash Tables

CSE 373
Data Structures & Algorithms
Ruth Anderson

2/11/2009

1

Today's Outline

- **Announcements**
 - Assignment #4 due this Friday Feb 13th at the beginning of lecture.
- **Today's Topics:**
 - **Disjoint Sets & Dynamic Equivalence**
 - **Hashing**

2/11/2009

2

Dictionary Implementations

	Unsorted linked list	Sorted Array	Binary Search Tree	AVL Tree
Insert				$O(\log N)$
Find	$O(N)$			
Delete			$O(N)$	$O(\log N)$

2/11/2009

3

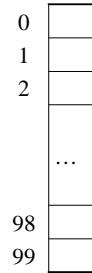
Constant Time Access

Data Set:

- 100 students
- Keys = Student numbers between 0 and 99.

Solution:

- Array of size 0-99.
- One-to-one mapping: e.g. student number 2 goes in location 2



2/11/2009

4

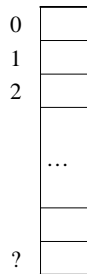
Constant Time Access?

Data Set:

- 100 students
- Keys = Student numbers between 0 and 999999999.

Solution:

- Array of size ?
- Mapping ?

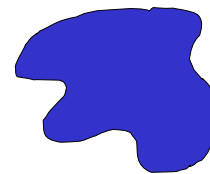


2/11/2009

5

Hash Tables

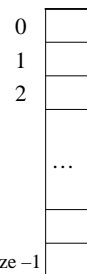
- A **hash table** is an array of some fixed size.
- General idea:



Hash Function:
 $h(K)$



Hash Table



Key Space (e.g., integers, strings)

TableSize - 1

2/11/2009

6

Example

- Key space = integers
- TableSize = 10
- $h(K) = K \bmod 10$
- **Insert:** 207, 18, 41, 194, 19, 43

0	
1	
2	
3	
4	
5	
6	
7	
8	
9	

2/11/2009 7

Another Example

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

0	
1	
2	
3	
4	
5	

2/11/2009 8

Student Activity

Hash Functions

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

Perfect Hash function:

2/11/2009 9

Sample Hash Functions:

key space = strings A=0, B=1,...Z=25

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

1. $h(s) = s_0 \bmod \text{TableSize}$
2. $h(s) = \left(\sum_{i=0}^{k-1} s_i \right) \bmod \text{TableSize}$
3. $h(s) = \left(\sum_{i=0}^{k-1} s_i \cdot 26^i \right) \bmod \text{TableSize}$

2/11/2009 10

Designing a Hash Function for web URLs

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

Issues to take into account:

$h(s) =$

2/11/2009 11

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)

2/11/2009 12

Separate Chaining

$h(K) = K \bmod 10$

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 (or "bucket").

2/11/2009 13

Analysis of Find

The **load factor**, λ , of a hash table is the ratio:

$$\frac{N}{\text{TableSize}} \leftarrow \begin{array}{l} \# \text{ of elements} \\ \text{table size} \end{array}$$

For separate chaining,
 $\lambda =$ average # of elements in a bucket

Average # of values needed to examine for a:

- unsuccessful find:
- successful find:

2/11/2009 14

How Big Should the Hash Table Be?

For Separate Chaining, if we want $\lambda = 1$
 (e.g. the average # of values per bucket = 1)

- How large should I make the hash table, in terms of N?

TableSize =

2/11/2009 15

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 ☺

2/11/2009 16

Open Addressing

$h(K) = K \bmod 10$

0	
1	
2	
3	
4	
5	
6	
7	
8	
9	

Insert:
38
19
8
109
10

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

2/11/2009 17

Terminology Alert!

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

2/11/2009 18

Linear Probing

$$f(i) = i$$

- Probe sequence:

$$0^{\text{th}} \text{ probe} = h(k) \bmod \text{TableSize}$$

$$1^{\text{th}} \text{ probe} = (h(k) + 1) \bmod \text{TableSize}$$

$$2^{\text{th}} \text{ probe} = (h(k) + 2) \bmod \text{TableSize}$$

...

$$i^{\text{th}} \text{ probe} = (h(k) + i) \bmod \text{TableSize}$$

2/11/2009

19

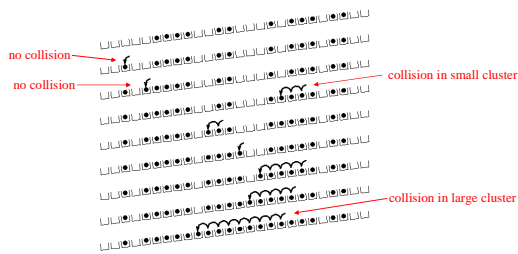
Write pseudocode for find(k) for Open Addressing with linear probing

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

Student Activity

20

Linear Probing – Clustering



[R. Sedgewick]

2/11/2009

21

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 $\lambda > 1/2$

2/11/2009

22

Quadratic Probing

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

Less likely
to encounter
Primary
Clustering

- Probe sequence:

$$0^{\text{th}} \text{ probe} = h(k) \bmod \text{TableSize}$$

$$1^{\text{th}} \text{ probe} = (h(k) + 1) \bmod \text{TableSize}$$

$$2^{\text{th}} \text{ probe} = (h(k) + 4) \bmod \text{TableSize}$$

$$3^{\text{th}} \text{ probe} = (h(k) + 9) \bmod \text{TableSize}$$

...

$$i^{\text{th}} \text{ probe} = (h(k) + i^2) \bmod \text{TableSize}$$

2/11/2009

23

Quadratic Probing

0		Insert:
1		89
2		18
3		49
4		58
5		79
6		
7		
8		
9		

2/11/2009

24

Quadratic Probing Example

insert(76) $76\%7 = 6$ insert(40) $40\%7 = 5$ insert(48) $48\%7 = 6$ insert(5) $5\%7 = 5$ insert(55) $55\%7 = 6$
 But... insert(47) $47\%7 = 5$

0	
1	
2	
3	
4	
5	
6	76

2/11/2009 25

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 $(i+j)$

2/11/2009 26

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!**

2/11/2009 27

Double Hashing

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

where g is a second hash function

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

2/11/2009 28

Double Hashing Example

$h(k) = k \bmod 7$ 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

2/11/2009 29

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

2/11/2009 30

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?

2/11/2009

31

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.

2/11/2009

32