

CSE 326: Data Structures

Hash Tables

Hal Perkins
 Spring 2007
 Lecture 16

1

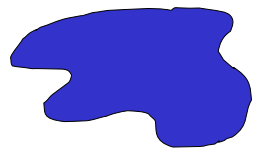
Dictionary Implementations So Far

	Unsorted linked list	Sorted Array	BST	AVL	Splay (amortized)
Insert					
Find					
Delete					

2

Hash Tables

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

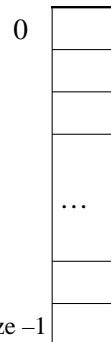


key space (e.g., integers, strings)

hash function:
 $h(K)$



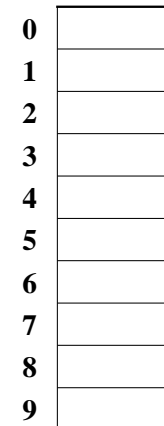
hash table



3

Example

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



4

Another Example

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

0	
1	
2	
3	
4	
5	

5

Hash Functions

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

Perfect Hash function:

6

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. $h(s) = \left(\sum_{i=0}^{k-1} s_i \cdot 37^i \right) \bmod \text{TableSize}$

7

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)

8

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 (or “bucket”).

9

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 find:
- Successful find:

10

How big should the hash table be?

- For Separate Chaining:

11

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 ☺

12

Open Addressing

0	
1	
2	
3	
4	
5	
6	
7	
8	
9	

Insert:
38
19
8
109
10

- **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.

13

Terminology Alert!

“Open Hashing” equals “Closed Hashing”

Weiss “Separate Chaining” “Open Addressing”

14

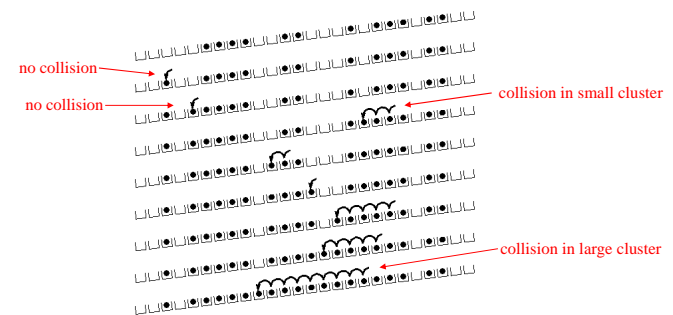
Linear Probing

$$f(i) = i$$

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

15

Linear Probing – Clustering



[R. Sedgwick]

16

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 $\lambda > 1/2$

17

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

18

Quadratic Probing

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

19

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)
 $47\%7 = 5$

20

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 $(i+j)$

21

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!*

22

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
$$(h(x) + i^2) \bmod \text{HSize} = (h(x) + j^2) \bmod \text{HSize}$$
$$(h(x) + i^2) - (h(x) + j^2) \bmod \text{HSize} = 0$$
$$(i^2 - j^2) \bmod \text{HSize} = 0$$
$$(i-j)(i+j) \bmod \text{HSize} = 0$$
$$\Rightarrow \Leftarrow \text{HSize does not divide } (i-j) \text{ or } (i+j)$$

23

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 = $(h(k) + 3 * g(k)) \bmod \text{TableSize}$
 - ...
 - i^{th} probe = $(h(k) + i * g(k)) \bmod \text{TableSize}$

24

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

25

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

26

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?

27

Java hashCode() Method

- Class Object defines a hashCode method
 - Intent: returns a suitable hashcode for the object
 - Result is arbitrary int; must scale to fit a hash table (e.g. `obj.hashCode() % nBuckets`)
 - Used by collection classes like HashMap
- Classes should override with calculation appropriate for instances of the class
 - Calculation should involve semantically “significant” fields of objects

28

hashCode() and equals()

- To work right, particularly with collection classes like HashMap, hashCode() and equals() must obey this rule:
 - if a.equals(b) then it must be true that
a.hashCode() == b.hashCode()
 - Why?
- Reverse is not required

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.

30