

VE 281 Homework 2 (Due 10/31 11:59PM)

Name: 赵海云 ID: 51837096091

• Exercise 1. Open Addressing

Consider inserting the keys 10, 22, 31, 4, 15, 28, 17, 88, 59 into a hash table of length $m = 11$ using open addressing with the auxiliary hash function $h(k) = k$. Illustrate the result of inserting these keys using linear probing, using quadratic probing with $c_1 = 1$ and $c_2 = 3$, and using double hashing with $h_1(k) = k$ and $h_2(k) = 1 + (k \bmod (m-1))$.

Solution:

We use T_t to represent each time stamp t starting with $i = 0$, and if encountering a collision, then we iterate i from $i = 1$ to $i = m - 1 = 10$ until there is no collision.

Linear probing:

$h(k, i) = (k + i) \bmod 11$	T_0	T_1	T_2	T_3	T_4	T_5	T_6	T_7	T_8
0 mod 11	22	22	22	22	22	22	22	22	22
1 mod 11							28	28	
2 mod 11									
3 mod 11									
4 mod 11	4	4	4	4	4	4			
5 mod 11		15	15	15	15	15	13		
6 mod 11			28	28	28	28			
7 mod 11				17	17	17			
8 mod 11							59		
9 mod 11	31	31	31	31	31	31	31	31	
10 mod 11	10	10	10	10	10	10	10	10	10

Quadratic probing:

$h(k, i) = (k + i + 3i^2) \bmod 11$	T_0	T_1	T_2	T_3	T_4	T_5	T_6	T_7	T_8
0 mod 11	22	22	22	22	22	22	22	22	22
1 mod 11							28	28	
2 mod 11									
3 mod 11						17	17	17	
4 mod 11	4	4	4	4	4	4	4	4	
5 mod 11									
6 mod 11			28	28	28	28			
7 mod 11				15	15	15	15	15	
8 mod 11					12	12	12	12	
9 mod 11	31	31	31	31	31	31	31	31	
10 mod 11	10	10	10	10	10	10	10	10	10

Double hashing:

$h(k, i) = (k + i(1 + k \bmod 10)) \bmod 11$	T_0	T_1	T_2	T_3	T_4	T_5	T_6	T_7	T_8
0 mod 11	22	22	22	22	22	22	22	22	22
1 mod 11									
2 mod 11								59	
3 mod 11						17	17	17	
4 mod 11	4	4	4	4	4	4	4	4	
5 mod 11			15	15	15	15	15	15	
6 mod 11				28	28	28	28	28	
7 mod 11						39	39	39	
8 mod 11									
9 mod 11	31	31	31	31	31	31	31	31	
10 mod 11	10	10	10	10	10	10	10	10	10

• Exercise 2. Slot-Size Bound for Chaining

Suppose that we have a hash table with n slots, with collisions resolved by chaining, and suppose that n keys are inserted into the table. Each key is equally likely to be hashed to each slot. Let M be the maximum number of keys in any slot after all the keys have been inserted. Your mission is to prove an $O(\lg n / \lg \lg n)$ upper bound on $E[M]$, the expected value of M .

a. Argue that the probability Q_k that exactly k keys hash to a particular slot is given by

$$Q_k = \left(\frac{1}{n}\right)^k \left(1 - \frac{1}{n}\right)^{n-k} \binom{n}{k}$$

b. Let P_k be the probability that $M = k$, that is, the probability that the slot containing the most keys contains k keys. Show that $P_k \leq nQ_k$.

c. Use Stirling's approximation to show that $Q_k < e^k/k^k$.

d. Show that there exists a constant $c > 1$ such that $Q_{k_0} < 1/n^3$ for $k_0 = c \lg n / \lg \lg n$. Conclude that $P_k < 1/n^2$ for $k \geq k_0 = c \lg n / \lg \lg n$

e. Argue that

$$E[M] \leq \Pr\left\{M > \frac{c \lg n}{\lg \lg n}\right\} \cdot n + \Pr\left\{M \leq \frac{c \lg n}{\lg \lg n}\right\} \cdot \frac{c \lg n}{\lg \lg n}$$

Conclude that $E[M] = O(\lg n / \lg \lg n)$

Solution:

a. $(\frac{1}{n})^k$ possibility of k keys inserted into a particular slot

$(1 - \frac{1}{n})^{n-k}$ all other keys not in this slot

$\binom{n}{k}$ choose k keys from n

$$\text{so } Q_k = \left(\frac{1}{n}\right)^k (1 - \frac{1}{n})^{n-k} \binom{n}{k}$$

b. nQ_k : there's one slot containing k keys (but may not be the only one)

$nQ_k = P[\text{one slot contains } k \text{ keys and all others } \overset{\text{no longer}}{\cancel{k+1 \text{ keys}}} \dots]$ + $P[\text{one slot contains } k \text{ keys and there are some more than } k]$

$$= P_k + P_{\neq k}$$

$$\text{so } P_k \leq nQ_k$$

$$C. \binom{n}{k} = \frac{n!}{k!(n-k)!} = \frac{\left(\frac{n}{e}\right)^n}{\sqrt{2\pi n} \left(\frac{k}{e}\right)^k \left(\frac{n-k}{e}\right)^{n-k}} = \frac{n^n}{\sqrt{2\pi n} k^k (n-k)^{n-k}}$$

$$Q_k = \frac{n^n}{\sqrt{2\pi n} k^k (n-k)^{n-k}} \cdot \frac{1}{n}^k \left(1 - \frac{1}{n}\right)^{n-k} = \frac{1}{\sqrt{2\pi n} k^k} \cdot \left(1 + \frac{k}{n-k}\right)^{-n+k} \cdot \left(1 - \frac{1}{n}\right)^{n-k} < \underbrace{\frac{1}{\sqrt{2\pi n}}}_{< 1} \cdot \underbrace{\left(\frac{e}{1-n}\right)^{n-k}}_{< 1} \cdot \frac{e^k}{k^k} < \frac{e^k}{k^k}$$

d.

$$\ln Q_{k_0} < -k \lg k \quad \ln Q_k < -k \lg k$$

$$\ln Q_{k_0} < -c(\ln n / \ln n) \cdot (\ln n - h \ln n)$$

$$\ln \ln n \gg \ln \ln n \quad \text{for large } n$$

$$\text{so } \ln Q_{k_0} < \frac{c}{2} \ln n$$

$$c=6 \Rightarrow Q_{k_0} < 3^{\ln n} \Rightarrow Q_{k_0} < \frac{1}{n^2}$$

$$P_{k_0} < n Q_{k_0} < \frac{1}{n^2}$$

e.

$$M_{\max} = n$$

$$E[M] = \sum \Pr\{M=i>0\} \cdot i + \sum \Pr\{M=i<0\} \cdot i$$

$$\leq \Pr\{M>0\} \cdot n + \Pr\{M<0\} \cdot 0$$

let $O = \frac{c \lg n}{\lg \lg n}$

~~For $M \geq \lg \lg n$~~

$$\leq \left(\sum_{i=0}^{n-1} \frac{1}{n^i} \cdot n \right) + \left(1 - \sum_{i=0}^{n-1} \frac{1}{n^i} \right) \cdot 0$$

$$\downarrow \qquad \qquad \downarrow$$

$$O(1) \qquad O(0)$$

$$\text{so } E[M] = O\left(\frac{c \lg n}{\lg \lg n}\right)$$

• Exercise 3. D-select

In the deterministic linear-time selection algorithm, instead of partitioning all the inputs into a number of groups of size 5, we partition the inputs into a number of groups of size 7. The rest steps of the algorithm are similar to the original version. Will the runtime of this new algorithm still be $O(n)$, where n is the input size? Prove your claim.

Solution:

Suppose there's a C such that

$$T_n \leq cn + T\left(\frac{n}{7}\right) + T\left(\frac{5n}{7}\right) \quad \leftarrow 1 - \frac{4}{7} \times \frac{1}{2}$$

and $T(1) \leq C$

We hope there's a constant a such that

$$T(n) \leq an \text{ for all } n > 1$$

We choose $a = 10C$ and then prove it

First, because $T(1) \leq C$, so $T(1) \leq 10C$

If $T(m) \leq am$ is true for all $m \leq k$

$$\begin{aligned} \text{then } T(k+1) &\leq c(k+1) + (k+1)c + 5(k+1)c \\ &= 7c(k+1) \end{aligned}$$

By induction, $T(n) \leq an$ for all pos int n (a.s.c.)

$$\text{therefore } T(n) = O(n)$$

• Exercise 4. Silly Blue Tiger

- Imagine that an algorithm requires us to hash strings containing English phrases. Knowing that strings are stored as sequences of characters, Blue Tiger decides to simply use the sum of those character values (modulo the size of its hash table) as the string's hash. Will the performance of the implementation match the expected value shown in lecture? (Yes/No and the brief reasoning are what we expect.)
- Blue Tiger decides to implement both collision resolution and rehashing for its hash table. However, it does not want to do more work than necessary, so it wonders if it needs both to maintain the correctness and performance it expects. After all, if it has rehashing, it can resize to avoid collisions; and if it has collision resolution, collisions don't cause correctness issues. Which statement about these two properties true? (Choose one or two of the statements and briefly explain your choice.)
- Blue Tiger wants to implement rehashing. Assuming it is enlarging a table of size m into a table of size m' , and the table contains n elements, what is the best time complexity to achieve this rehashing step (expanding a m -slot table with n elements into a m' -slot table)? Answer in big theta notation in terms of the variables in the context. Intermediate results of your complexity analysis should be displayed.
- In lecture, we discussed approximately doubling the size of our hash table. Blue Tiger begins to implement this approach (that is, it lets m' be a prime close to but greater than $2m$) but stops when it thinks that it might be able to avoid wasting half of the memory the table occupies on empty space by letting m' be a prime number close to $(m + k)$ instead, where k is some constant. Does this work equally well comparing to picking a prime number close to $2m$? Explain your answer.

Solution:

1. no. words with same characters but different sequence will be hashed to same slot "top" "pot"
2. collision resolution ~~will~~ let collision due to cause correctness issue is there.
It will just take longer time (longer and longer) to ~~be~~ carry out.
rehashing reduce the possibility of collision but cannot avoid it

?

• Exercise 5. Rehashing

6/26 - 02/7 - 96/2 - 14/2

Given a sequence of inputs 4371, 1323, 6173, 4199, 4344, 9679, 1989, insert them into a hash table of size 10.

Suppose that the hash function is $h(x) = x \bmod 10$. Hash table uses the double hashing with the second hash function as $g(x) = (7 - x) \bmod 7$. (Be careful with this mod calculation) $[h(x) + i g(x)] \bmod n$

If any given input cannot be inserted correctly, do rehashing to guarantee a successful insertion, using a new table size as introduced in class (the first prime number after doubling the original hash table size) and a new hash function $h(x) = x \bmod \text{new_size}$. Note that the order in rehashing depends on the order stored in the old hash table, not their previous inserting order.

Solve this question step by step and show intermediate results.

Solution:

The initial hash table:

	T_0	T_1	T_2	T_3	T_4	T_5	T_6
entry0							
entry1	4371	4371	4371	4371	4371	4371	
entry2							
entry3		1323	1323	1323	1323	1323	
entry4				6173	6173	6173	
entry5						9679	
entry6							
entry7					4199	4199	
entry8							1989

rehash

Does the hash table need rehashing? If it does, when and why (answer in 1-2 sentences)? Also finish the rehashing process by drawing your own progress table that is similar to the one above:

$$T_6 \quad h_i = (9 + 6i) \% 10, \quad 13579 \text{ are occupied}$$

	T_0	T_1	T_2	T_3	T_4	T_5	T_6
0							
1	4371						
2							
3							
4							
5							
6							
7							
8							
9			6173				
10							
11							
12		1323					
13							
14							
15							
16							
17							
18							
19						9679	
20							4199
21							
22							
23							

• Exercise 6

Suppose we want to design a hash table containing at most 1162 elements using quadratic probing. We require it to guarantee successful insertions, $S(L) < 2$ and $U(L) < 4$. Please determine a proper hash table size and show all intermediate steps.

$$\begin{aligned} \frac{1}{2} \left(\frac{L+1}{L} \right)^2 L &< 2 \\ \frac{1}{2} \left(\frac{14}{12} \right)^2 L &< 2 \Rightarrow L < 6.62 \quad \text{so } L < 7.25 \\ \frac{1}{L} < U &\Rightarrow L < 0.75 \\ \left\{ \frac{1}{L} \ln \frac{1}{1-L} < 2 \right\} &\Rightarrow L < 0.79 \quad \text{take the closest prime } 11 \end{aligned}$$

• Exercise 7

In a few sentences (in an itemized fashion), explain why a hash table with open addressing and rehashing (doubling the size every time when the load factor exceeds a threshold) achieves insertion in $O(1)$ time.

Answer the following questions: 1) When rehashing is not triggered, is the insertion time $O(1)$ (does it have an upper bound)? Why? 2) Show that at any given point, the amortized (distributed to all items in the hash table) cost of rehashing is $O(1)$ for each item in the hash table.

Because open addressing cost $O(1)$ time, and rehashing cost $O(4)$ time amortize over each insert.

1. hashing: $O(1)$
2. open addressing: $O(1)$ before rehashing
3. rehash: $O(1)$ if amortized.

- (1) No. the worst case is $O(n)$ e.g. linear probing, and all elements have the same bucket original cost
- (2) suppose size $M \rightarrow 2M$, and we have ML items already. (L load factor). Adding another triggers rehash, we need to insert $(M+1)$ items, cost $O(M)$ time. Rate $\frac{2O(M)}{LM+1} = O(1)$

What if during rehashing, the hash table is expanded by a factor x with $x > 2$? Is the amortized cost of rehashing still $O(1)$? If no, explain. If yes, show your work. If it depends, show both.

Yes. suppose each time the table expanded x^X times, & suppose the last rehashed table is of size M . then the total time on rehashing is $O(M) + \frac{1}{x} O(M) + \frac{1}{x^2} O(M) - \dots$
 $= O(M)$, and the hashed items number is $O(M)$
 so still $O(1)$, if amortized.