

Data Structures for Dictionaries: Hash Tables

Problem 1: Read a grade file, where grades are integers between 0 to 99. Keep track of number of occurrences of each grade.

Fastest Solution: Create an array T of size 100. $T[i]$ stores the number of occurrences of grade i .

Problem 2: Read a data file, keep track of number of occurrences of each integer value (from 0 to $2^{32} - 1$).

Fastest Solution: Create an array of size 2^{32} , as above.

Wasteful use of memory, especially when data are files relatively small.

Problem 3: Read a text file, keep track of number of occurrences of each word.

Cannot use **keys as indices** anymore!

1. We need to be able to convert any type of key to an integer.
2. We need to map the universe of keys into a small number of slots.

csc263 lec4 ppt pg3

研究：

- 增 (插入 insert)
- 删 (delete)
- 查 (寻址/search)

List 是 ADT ADT，其两种基本的实现方式有 array 数组， linkedlist 链表

array[4]=60, 4 是 index 索引， 60 是其对应的数据、数值、value
index 索引是提前设置好的

HashTable 基本原理包含了 array，随机顺序访问数据/查询，**访问查询速度最快**，快速根据 key(index 索引) $O(1)$

Array 是简单的 HashTable

HashTable: 在直接寻址下，具有关键字 k 的元素被存放在槽 k 中。在散列方式下，该元素存放在 $h(k)$ 中；即利用散列函数 (hash function) h ，由关键字 k 计算出槽的位置。—教材 p143

LinkedList 链式一般是从第一个进行 traverse, $O(n)$

数据源 (file)– 存储(array)

将数据源里的 type of key 转换成 integer (通过 hash function 转换成具有整数性质的散列值)，引出 Hash Function (universe of keys)

- slots, slot 里面即每个槽位里存的还是原关键字 key)

Hash Table (散列表/哈希表) : 元素, 关键字 (key), 散列值 (h(k))

承载因子 (load factor) 一般是 0.75

目标：使数据尽可能均匀放在 slots 里，ensure uniform

Expected Run Time in a Successful Search (under SUHA):

That is k is a key that exists in the hash table.

Let $k_1, k_2, k_3, \dots, k_n$ be the order of insertion into the hash table.

k could be k_1 , or k_2 , or k_3 , or ..., or k_n .

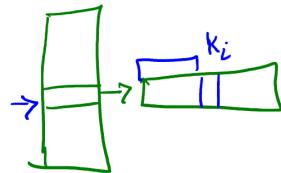
The probability that k is k_i ($1 \leq i \leq n$) is: $\frac{1}{n}$

So the expected number of steps to find k is the sum over:

the probability that k is k_i , times the number of steps required to find k_i

$$\begin{aligned}\mathbb{E}[t_{m,n}(k)] &= \frac{1}{n} \times S_1 + \frac{1}{n} \times S_2 + \frac{1}{n} \times S_3 + \dots + \frac{1}{n} \times S_n \\ &= \frac{1}{n} \sum_{i=1}^n S_i\end{aligned}$$

S_i denotes the expected number of steps to find k_i .



S_i = expected number of steps to find k_i

. = number of elements examined during search for k_i

. = $1 +$ number of elements before k_i in the linked list stored at $h(k_i)$

. = $1 +$ number of keys that hash similarly as k_i and are inserted after k_i .

csc263 lec4 ppt pg14

注意：链接法（chaining）中的插入和搜索的顺序（不同）

“在对元素 x 的一次成功查找中，所检查的元素就是 x 所在的链表中 x 前面的元素数多 1。在该链表中因为，新的元素都是在表头插入的，所以出现在 x 之前的元素都是 x 之后插入的” - 教材 p146

chaining search avg - successful case,

$$\begin{aligned}\mathbb{E}[t_{m,n}(k)] &= \frac{1}{n} \sum_{i=1}^n S_i \\ &= \frac{1}{n} \left[1 + \sum_{j=i+1}^n \mathbb{E}[X_{i,j}] \right] \\ &= \frac{1}{n} \left[1 + \sum_{j=i+1}^n \left(\frac{1}{m} \right) \right] \xrightarrow{\text{因为在 slot } j \text{ 有 } k_i} \\ &= \frac{1}{n} \left[n + \sum_{i=1}^{n-1} \frac{n-i+1}{m} \right] \xrightarrow{\text{其 slot 分布 probability}} \\ &= \frac{1}{n} \left[n + \cancel{\frac{n(n+1)}{2m}} \right] \xrightarrow{\cancel{\frac{n(n+1)}{2m}}} \frac{n \cdot n - \frac{(n+1)n}{2m}}{2m} \\ &= 1 + \frac{n}{m} - \frac{1}{m} - \frac{(n+1)n}{2m} \\ &= 1 + \frac{n}{2m} - \frac{1}{2m} \in \Theta(n^2).\end{aligned}$$

Search 的方向 $\rightarrow \rightarrow \rightarrow$
但并非是 insert 的方向

insert 的方向 $\leftarrow \leftarrow \leftarrow$

Diagram illustrating the search and insertion directions in a linked list. The top part shows a search path starting from the head of a list and moving right. The bottom part shows an insertion operation where a new node k_1 is inserted at index i before node k_i . The list is labeled as being implemented as a linked list at the end of the array.

定理 11.1 在 SUHA (简单均匀散列) 的假设下, 对于用链接法解决冲突的散列表, 一次不成功查找的平均时间为 big theta($1+\alpha$)

定理 11.2 在 SUHA (简单均匀散列) 的假设下, 对于用链接法解决冲突的散列表, 一次成功查找的平均时间为 big theta($1+\alpha$)

Hash Function 散列函数/哈希函数 与 Probing 定位

Open Addressing: Linear probing

Linear Probing: Examine a linear sequence of slots.

Example probe sequence: $h(k), h(k) + 1, h(k) + 2, h(k) + 3, \dots$

Probe sequence: $(h(k) + i) \bmod m$, for an integer $i \geq 0$.

Open Addressing: Quadratic Probing

Quadratic Probing: Examine a non-linear sequence of slots.

Example probe sequence: $h(k), h(k) + 2, h(k) + 6, h(k) + 12, \dots$

Probe sequence: $(h(k) + c_1 \times i + c_2 \times i^2) \bmod m$, for an integer $i \geq 0$.

Open Addressing: Linear probing

Linear Probing: Examine a linear sequence of slots.

Example probe sequence: $h(k), h(k) + 1, h(k) + 2, h(k) + 3, \dots$

Probe sequence: $(h(k) + i) \bmod m$, for an integer $i \geq 0$.

ppt lecture4 pg22-24

注意这里的 example probe sequence 并不是但单纯的 example, 而是常态下的基本 probing 序列

TUT 例题

2. In this question, you will explore detail implementation for hash table operations when the table uses open addressing.
 - (a) Insert the keys 10, 22, 31, 4, 15, 28, 17, 88, 39 into a hash table of size $m = 11$ using open addressing with the primary hash function $h_1(k) = (k \bmod m)$. Illustrate the results for linear probing, quadratic probing, and double hashing with $h_2(k) = 1 + (k \bmod (m - 1))$.

primary hash function: $h_1(k) = (k \bmod m)$ $m=11$

list index 0 1 2 3 4 5 6 7 8 9 10
 h₁(k) 22 88 4 15 28 17 39 31 10

h₁(k) 使用 primary hash function.

using linear probing: probe sequence {h₁(k), h₁(k)+1, h₁(k)+2, ...}

$$10 \bmod 11 = 10, 22 \bmod 11 = 0, 31 \bmod 11 = 9$$

$$4 \bmod 11 = 4, 15 \bmod 11 = 4, 28 \bmod 11 = 6$$

$$17 \bmod 11 = 6 \rightarrow \text{collision} \rightarrow 6+1=7 \rightarrow \text{collision} \rightarrow 4+1=5$$

$$88 \bmod 11 = 0 \rightarrow \text{collision} \rightarrow 0+1=1$$

$$39 \bmod 11 = 6 \rightarrow \text{collision} \rightarrow 6+1=7 \rightarrow \text{collision} \rightarrow 7+1=8$$

quadratic probing: probe sequence {h₁(k), h₁(k)+1, h₁(k)+2, h₁(k)+3, ...}

h₁(k) 使用 primary hash function.

list index 22 88 4 15 28 17 39 31 10

$$15 \bmod 11 = 4 \rightarrow \text{collision} \rightarrow 4+1^2=5$$

$$17 \bmod 11 = 6 \rightarrow \text{collision} \rightarrow 6+1^2=7$$

$$88 \bmod 11 = 0 \rightarrow \text{collision} \rightarrow 0+1^2=1$$

$$29 \bmod 11 = 6 \rightarrow \text{collision} \rightarrow 6+1^2=7 \rightarrow \text{collision} \rightarrow 6+2^2=10 \bmod 11 = 1$$

$$\text{collision} \leftarrow 6+4^2 \bmod 11 = 0 \quad \xleftarrow{\text{collision}} \quad 6+3^2 \bmod 11 = 4 \quad \xleftarrow{\text{collision}}$$

$$6+5^2 \bmod 11 = 9 \rightarrow \text{collision} \rightarrow 6+6^2 \bmod 11 = 9 \rightarrow \text{collision} \rightarrow 6+7^2 \bmod 11 = 0$$

$$6+10^2 \bmod 11 = 7 \quad \leftarrow \quad 6+9^2 \bmod 11 = 10 \quad \leftarrow \quad 6+8^2 \bmod 11 = 6 \quad \leftarrow \quad \text{collision}$$

although not full in list, we have to stop to probe b/c always collision and we can't insert...

h₁(k) double hashing using double hashing if h₁(k) fails

list index 0 1 2 3 4 5 6 7 8 9 10
 h₁(k) 22 17 4 15 28 88 31 10

$$\text{hw: } 15 \bmod 11 = 4 \rightarrow \text{collision} \quad 4 + h_2(15) = 4 + 6 = 10 \rightarrow \text{collision}.$$

$$17 \bmod 11 = 6 \rightarrow \text{collision} \quad 6 + h_2(17) \quad \swarrow +2h_2(15)$$

$$= 6 + 8 = 14 \text{ index} \quad 10 + 6 = 16 \text{ index}$$

从6起向右
数8个 index

$$88 \bmod 11 = 0 \rightarrow \text{collision} \quad 0 + h_2(88) = 0 + 9 = 9 \rightarrow \text{collision}$$

$$0 + 2h_2(88) = 0 + 18 = 18 \text{ index}$$