

# 论题 2-10 作业

姓名：陈劭源

学号：161240004

## 1 [CS] Problem 5.8

a. When  $n \leq k$ , the probability that all  $n$  items hash to different locations is

$$P_1 = \frac{k^n}{k^n} = \frac{k!}{(k-n)!k^n}$$

and when  $n > k$ , the probability is 0.

b. When  $i \leq k+1$ , the probability that the  $i$ th item is the first collision is

$$P_2(i) = \frac{k^{i-1}}{k^{i-1}} \times \frac{i-1}{k} = \frac{k!(i-1)}{(k-i+1)!k^i}$$

and when  $i > k+1$ , the probability is 0.

c. Let  $X$  denote the number of items you hash until the first collision, then

$$E(X) = \sum_{i=2}^{k+1} (i-1)P_2(i) = \left( \sum_{i=2}^{k+1} \frac{k!(i-1)i}{(k-i+1)!k^i} \right) - 1$$

d. When  $k = 20$ , the expected number is 5.29358.

When  $k = 100$ , the expected number is 12.20996.

## 2 [CS] Problem 5.11

a. Since we must go through the whole list to ensure that the item is indeed not in the hash table, the running time for an unsuccessful search is  $\Theta(1 + \text{length})$ , where the expectation of  $\text{length}$  is  $n/k$ , thus the expected running time is  $\Theta(1 + n/k)$ .

b. Assume there are  $i$  elements inserted after the item you are searching for, then the expected running time is

$$t_i = \Theta\left(1 + \sum_{j=0}^i j \binom{i}{j} (1/k)^j (1-1/k)^{i-j}\right) = \Theta(1 + i/k)$$

And the expected running time for a successful search is

$$\Theta\left(\frac{1}{n} \sum_{i=0}^{n-1} t_i\right) = \Theta\left(\frac{1}{n} \sum_{i=0}^{n-1} (1 + i/k)\right) = \Theta(1 + (n-1)/2k)$$

Since there is no need to go through the whole list and we can stop searching immediately we've found the item, the expected running time of a successful search is about half of the unsuccessful one.

### 3 [CS] Problem 5.14

By Theorem 5.15, the expected number of empty slots when hashing  $2k$  items into a hash table with  $k$  slots is

$$E(X) = k(1 - 1/k)^{2k}$$

and when  $k$  grows large, the expected fraction of empty slots is

$$\lim_{k \rightarrow \infty} E(X)/k = \lim_{k \rightarrow \infty} k(1 - 1/k)^{2k}/k = \lim_{k \rightarrow \infty} (1 - 1/k)^{2k} = e^{-2}$$

### 4 [TC] Problem 11.2-3

Since the linked list is not randomly accessible, binary search does not apply and we can only use sequential search, so the modification does not affect the running time for successful searches and unsuccessful searches.

We have to insert the item to the correct position in the list, so the average running time of insertions becomes  $\Theta(1 + \alpha/2)$  where  $\alpha$  is the load factor.

This modification does not affect the running time for deletions, which is still  $O(1)$  if the lists are doubly linked.

### 5 [TC] Problem 11.2-5

There are  $|U| > nm$  elements (pigeons), but there are only  $m$  slots (holes). By the pigeonhole principle, there exists a subset of  $U$ , consisting of  $n$  keys that all hash to the same slot. If we want to search for the last element in the slot, the running time is  $\Theta(n)$ . Therefore, the worst-case searching time for hashing with chaining is  $\Theta(n)$ .

### 6 [TC] Problem 11.2-6

RANDOM-SELECT( $T, m, L$ )

```
1  repeat
2       $x = \text{RANDOM}(1, m)$ 
3       $y = \text{RANDOM}(1, L)$ 
4  until  $y \leq T[x].\text{length}$ 
5   $t = T[x].\text{head}$ 
6  for  $i = 2$  to  $y$ 
7       $t = t.\text{next}$ 
8  return  $t$ 
```

In line 1-4, we repeat trying to choose a position in the hash table randomly, and every position is equally likely to be chosen, so the procedure chooses the key uniformly at random.

For each iteration in line 1-4, the probability of success is  $n/mL = \alpha/L$ . So the total number of iterations follows the geometric distribution with parameter  $p = \alpha/L$ , and its expectation is  $1/p = L/\alpha$ . In line 5-7, we

search the element in a list, whose length is at most  $L$ , so the running time is  $O(L)$ . Therefore, the total expected running time is  $O(L + L/\alpha) = O(L(1 + 1/\alpha))$ .

## 7 [TC] Problem 11.3-3

The string can be represented as

$$S = \sum_{i=0}^{n-1} S_i (2^p)^i$$

Consider  $(2^p)^i \bmod (2^p - 1)$ , we have

$$(2^p)^i \bmod (2^p - 1) = (2^p \bmod (2^p - 1))^i \bmod (2^p - 1) = 1^i \bmod (2^p - 1) = 1$$

therefore

$$S \bmod (2^p - 1) = \sum_{i=0}^{n-1} S_i (2^p)^i \bmod (2^p - 1) = \sum_{i=0}^{n-1} S_i \bmod (2^p - 1)$$

The result is only in terms of the sum of the characters in the string and has nothing to do with the position of each character. That means, if string  $y$  is a permutation of string  $x$ , then they hash to the same value.

If there are many string pairs that one of them is a permutation of the other, such as ‘POT’ and ‘TOP’, in our application, then the hash function will lead to a lot of collisions, which is undesirable.

## 8 [TC] Problem 11.3-4

$k$	$kA$	$kA \bmod 1$	$m(kA \bmod 1)$	$h(k)$
61	37.70007	.70007	700.07	700
62	38.31811	.31811	318.11	318
63	38.93614	.93614	936.14	936
64	39.55418	.55418	554.18	554
65	40.17221	.17221	172.21	172

## 9 [TC] Problem 11.4-2

HASH-DELETE( $T, k$ )

1  $T[\text{HASH-SEARCH}(T, k)] = \text{DELETED}$

HASH-INSERT( $T, k$ )

```

1   $i = 0$ 
2  repeat
3       $j = h(k, i)$ 
4      if  $T[j] == \text{NIL}$  or  $T[j] == \text{DELETED}$ 
5           $T[j] = k$ 
6          return  $j$ 
7      else
8           $i = i + 1$ 
9  until  $i == m$ 
10 error “hash table overflow”

```

## 10 [TC] Problem 11.4-3

By Theorem 11.6 and Theorem 11.8, when the load factor is  $3/4$ , the expected number of probes in an unsuccessful search is at most  $1/(1 - \alpha) = 4$ , the expected number in a successful search is at most  $(1/\alpha) \ln 1/(1 - \alpha) = 1.848$ ; when the load factor is  $7/8$ , the expected number of probes in an unsuccessful search is at most  $1/(1 - \alpha) = 8$ , the expected number in a successful search is at most  $(1/\alpha) \ln 1/(1 - \alpha) = 2.377$ .

## 11 [TC] Problem 11-1

- a.** If an insertion requires strictly more than  $k$  probes, the first  $k$  probes are unsuccessful. For the  $(i + 1)$ th probe, the probability of unsuccessful probe is  $(n - i)/(m - i) \leq 1/2$  because the probe sequence is a random permutation of  $0, 1, \dots, m - 1$  according to uniform hashing assumption, and the probability that all the  $k$  probes are unsuccessful is less than  $2^{-k}$ . So the probability is at most  $2^{-k}$ .
- b.** Substituting  $k = 2 \lg n$  into **a.** yields that the probability is at most  $2^{-2 \lg n} = n^{-2}$ , i.e. the probability is  $O(1/n^2)$ .
- c.** Consider the complementary event, that  $X \leq 2 \lg n$ , which is equivalent to the event that for all  $i$ ,  $X_i \leq 2 \lg n$ . Because  $X_1, X_2, \dots, X_n$  are independent, we have

$$\begin{aligned}
 P(X \leq 2 \lg n) &= \prod_{i=1}^n (1 - P(X_i > 2 \lg n)) \\
 &\geq (1 - n^{-2})^n \\
 &\geq 1 - 1/n \quad \text{(By Bernoulli's inequality)}
 \end{aligned}$$

Hence,

$$P(X > 2 \lg n) = 1 - P(X \leq 2 \lg n) \leq 1/n$$

Therefore,  $P(X > 2 \lg n) = O(1/n)$ .

- d.** Since the probe sequence is a permutation of  $1, 2, \dots, m - 1$ , the maximum length of the probe sequence is the items in the hash table plus one, which is at most  $n$ . Hence, the expected length of the longest probe is

$$\begin{aligned}
 E[X] &= \sum_{i=1}^n iP(X = i) \\
 &= \sum_{i=1}^{\lfloor 2 \lg n \rfloor} iP(X = i) + \sum_{i=\lfloor 2 \lg n \rfloor + 1}^n iP(X = i) \\
 &\leq \lfloor 2 \lg n \rfloor \sum_{i=1}^{\lfloor 2 \lg n \rfloor} P(X = i) + n \sum_{i=\lfloor 2 \lg n \rfloor + 1}^n P(X = i) \\
 &= \lfloor 2 \lg n \rfloor \sum_{i=1}^{\lfloor 2 \lg n \rfloor} P(X = i) + nP(X > 2 \lg n) \\
 &\leq 2 \lg n + 1
 \end{aligned}$$

Therefore,  $E[X] = O(\lg n)$ .

## 12 [TC] Problem 11-2

- a.** Let random variable  $X_i$  denote the number of keys hashing to the  $i$ th slot.  $X_i$  follows the binomial distribution with parameter  $n$  and  $p = 1/n$ . Therefore

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

- b.** Let  $E_{i,k}$  denote the event that the  $i$ th slot contains the most keys and it contains exactly  $k$  keys. Since  $E_{i,k} \subset (X_i = k)$ , we have  $P(E_{i,k}) \leq Q_k$ . Therefore,

$$P_k = P(M = k) = P\left(\bigcup_{i=1}^n E_{i,k}\right) \leq \sum_{i=1}^n P(E_{i,k}) \leq nQ_k$$

- c.** First, we have to bound the binomial coefficient

$$\begin{aligned} \binom{n}{k} &= \frac{n!}{(n-k)!k!} \\ &\leq \frac{n^k}{k!} \\ &\leq n^k \left(\frac{e}{k}\right)^k \end{aligned} \quad \text{(By equation (3.18), Stirling's approximation)}$$

Then we have

$$\begin{aligned} Q_k &= \left(\frac{1}{n}\right)^k \left(1 - \frac{1}{n}\right)^{n-k} \binom{n}{k} \\ &< \left(\frac{1}{n}\right)^k \left(1 - \frac{1}{n}\right)^{n-k} n^k \left(\frac{e}{k}\right)^k \\ &\leq \left(\frac{1}{n}\right)^k n^k \left(\frac{e}{k}\right)^k \\ &= e^k / k^k \end{aligned}$$

- d.** Assume that  $n \geq 3$ , otherwise  $c \lg n / \lg \lg n$  is undefined. Consider

$$\begin{aligned} \frac{\lg Q_{k_0}}{\lg n} &< \frac{k_0 \lg e - k_0 \lg k_0}{\lg n} \\ &= \frac{c \lg n (\lg e - \lg(c \lg n / \lg \lg n))}{\lg n \lg \lg n} \\ &= \frac{c(\lg e - \lg c - \lg \lg n + \lg \lg \lg n)}{\lg \lg n} \\ &= c \left( \frac{\lg e - \lg c}{\lg \lg n} - 1 + \frac{\lg \lg \lg n}{\lg \lg n} \right) \\ &\leq c \frac{\lg e - \lg c}{\lg \lg 3} \end{aligned}$$

Take  $c = 2e$ , we get  $\lg Q_{k_0} / \lg n < -4e / \lg \lg 3 < -3$ , i.e.  $Q_{k_0} < 1/n^3$ .

Since  $Q_{k+1}/Q_k = (n-k)/((n-1)(k+1)) \leq 1$  when  $k \geq 1$ ,  $\{Q_k\}$  is a monotonously decreasing sequence, and by **b.** we get  $P_k < 1/n^2$  for  $k \geq k_0 = c \lg n / \lg \lg n$ .

**e.** We have

$$\begin{aligned}
E[M] &= \sum_{k=1}^n kP_k \\
&= \sum_{k=1}^{\lfloor \frac{c \lg n}{\lg \lg n} \rfloor} kP_k + \sum_{k=\lfloor \frac{c \lg n}{\lg \lg n} \rfloor + 1}^n kP_k \\
&\leq P\left(M > \frac{c \lg n}{\lg \lg n}\right) n + P\left(M \leq \frac{c \lg n}{\lg \lg n}\right) \frac{c \lg n}{\lg \lg n}
\end{aligned}$$

and by **b.**, we have

$$E[M] < n \cdot 1/n^2 \cdot n + \frac{c \lg n}{\lg \lg n} = O(\lg n / \lg \lg n)$$