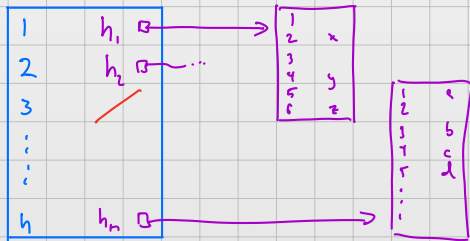


CS 5112

Streaming Algorithms

FKS Perfect Hashing

Level 1:



Level 2: Each bucket is a really big hash table (i.e. collision tree)

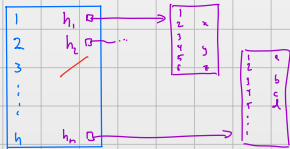
Let C_t be the # of items in bucket t .

The level 2 table for t has C_t^2 slots.

FKS Perfect Hashing

Level 1:

Level 2: Each bucket is a really big hash table (i.e. collision tree)



Let C_t be the # of items in bucket t .

The level 2 table for t has C_t^2 slots.

$O(n)$

Construct: 1. Loop through all items and hash them to buckets (lists)

Expected

$O(1)$

actually 2

2. For each bucket: a. Choose a random h_t . $O(1)$

$O(C_t)$

b. Try to assign items to level 2 slots using h_t .

c. If there's a collision, - back to a.

FKS Perfect Hashing

Level 1:

1	h_1	□
2	h_2	□
3		□
⋮		□
⋮		□
n	h_n	□

1	a
2	x
3	
4	y
5	
6	z

1	a
2	b
3	c
4	d
5	
⋮	

Level 2: Each bucket is a really big hash table (i.e. collision free)

Let C_t be the # of items in bucket t .

The level 2 table for t has C_t^2 slots.

expected

$$O(n) + O(1) \cdot \sum_{t=1}^n C_t$$

$$= O(n) + O(n) \text{ in expectation}$$

Construct: 1. Loop through all items and hash them to buckets (lists) $O(n)$

Expected $O(1)$ actually 2 { 2. For each bucket: a. Choose a random h_t $O(1)$
b. Try to assign items to level 2 slots using h_t $O(C_t)$
c. If there's a collision, go back to a.

FKS Perfect Hashing

What about space?

Level 1: Uses $\mathcal{O}(n)$

Level 2: $\sum_{t=1}^n C_t^2$ space ($2_{n+1} = \mathcal{O}(n)$).

$$E\{C_t^2\}$$

C_t = RV for # of items in level-1 bucket t .

$I_{i,t}$ = RV for whether $h(x_i) = t$. $\begin{cases} 1 & \text{if } h(x_i) = t \\ 0 & \text{otherwise} \end{cases}$

$$C_t = \sum_{i=1}^n I_{i,t}$$

$$E[C_t] = \sum_{i=1}^n E\{I_{i,t}\} = n/n = 1$$

Want something similar for C_t^2 .

$$J_{i,j,t} = \text{RV indicating whether } h(x_i) = t \wedge h(x_j) = t$$

$$= \begin{cases} 1 & \text{if } h(x_i) = t \text{ and } h(x_j) = t \\ 0 & \text{otherwise} \end{cases}$$

$$C_t^2 = \sum_{i,j} J_{i,j,t}$$

$$E\{C_t^2\} = \sum_{i,j} E\{J_{i,j,t}\}$$

$$E\left\{\sum_{t=1}^m C_t^2\right\} = \sum_{i,j} E\{J_{i,j,t}\} = \sum_{i,j} E\left\{\sum_{t=1}^m J_{i,j,t}\right\} = \sum_{i,j} \Pr(h(x_i) = h(x_j))$$

		a	b	c	d	...
i's	a	
	b	
	c	
	d	
	:					
		j's				

$$\begin{aligned}
E\left\{\sum_{\ell=1}^n C_{\ell}^2\right\} &= \sum_{(i,j) \in \mathcal{E}} E\{J_{i,j,1\ell}\} = \sum_{i,j} E\left\{\sum_{\ell=1}^n J_{i,j,\ell}\right\} = \sum_{i,j} \Pr(h(x_i) = h(x_j)) \\
&= \sum_{i=j} 1 + \sum_{i \neq j} \Pr(h(x_i) = h(x_j)) \\
&= n + \frac{1}{n} \sum_{i \neq j} 1 = n + (n^2 - n) \frac{1}{n} \\
&= 2n + 1.
\end{aligned}$$

Streaming Algorithms

Stream \rightarrow

1	2	3	...									
A	C	R	L	A	B	X	Z	D	A	A	A	B

The stream is really big, so we can't record it.

Still want to answer questions:

1. What's the most common item?
2. A list of frequent items
3. # unique items

Distinct Element Problem

First attempt: Use a little hash table:

A

B

A

Distinct Element Problem

Flajolet-Martin

Clever idea: Hash each item, ^{to $\{1, \dots, m\}$} keep track of the smallest hash, s .

If we know what the smallest hash is,
we can estimate the # of unique items.

Return c^m/s .

Problem: lots of variance.

Majority Problem

Input: A stream

Output: The majority element if it exists. (and whether one exists)

Problem: Can't determine if a majority exist with $< \frac{n}{2}$ space.

Boyer-Moore

On i th item:

Store: candidate item a
counter c

Majority Problem

Boyer-Moore

Store: candidate item a
counter c

On i th item x_i :

$a \neq \emptyset$
if $a = x_i$, $c++$
else $c--$
if $c = 0$, set $a = \emptyset$

$a = \emptyset$
set $a = x_i$
 $c++$



A B A C D F A A A A B C

Candidate a: ~~∅~~ ~~A~~ ~~∅~~ A ~~∅~~ ~~D~~ ~~∅~~ A

Counter c: ~~∅~~ X ~~∅~~ X ~~∅~~ X ~~∅~~ X ~~∅~~ X ≠ 3

Boyer-Moore

Suppose there is a majority element, B .

Case 1: C never goes back to 0. \rightarrow we see first item more
than all other items
combined

Case 2: C does go back to 0.

\Rightarrow first item is B .

\hookrightarrow Let's look at the last time $t_n^t C$ goes to 0.

1 2 . . . t | $t+1$. . . n

return the
majority item
from
 $\{t+1, \dots, n\}$