

1	2	3	10	11	12	13	14
4	5	6	17	18	19	20	21
7	8	9	23	24	25	26	27
10	11	12	28				
13	14	15					
16	17	18					
19	20	21					
22	23	24					
25	26	27					
28	29	30					

WR 40 • 275-090

02

OCTOBER

SATURDAY

Reservoir Sampling

- How to pick a random elements of an array?
- Ans: choose a random integer and pick the element at that index.
- How to pick 'S' random elements from an infinite stream of elements?
- Maintain a set of 'S' elements that may be updated every time a new element is added to the stream.
- At any given point, each element that has been seen so far should be equally likely to be in the set of S elements.

↓ How?

- Save the first S elements

03

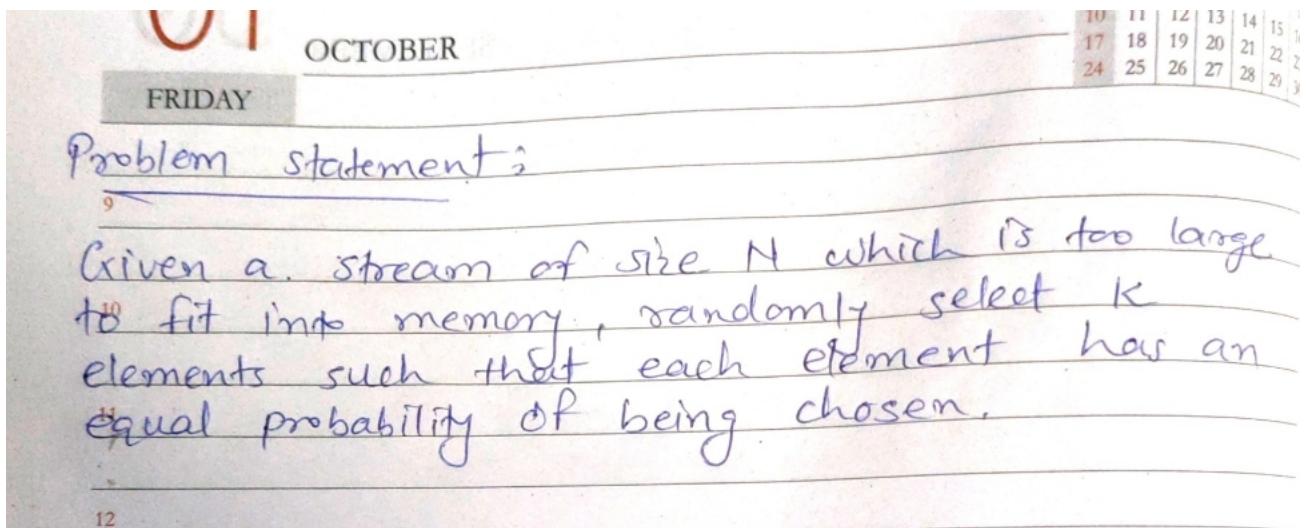
- Suppose we have seen the first  $n-1$  elements and now  $n^{\text{th}}$  element arrives.

- with probability  $S/n$ , keep  $n^{\text{th}}$  element and replace a randomly chosen saved element with the  $n^{\text{th}}$  element.
- otherwise discard the  $n^{\text{th}}$  element

2021 Induction :- After  $n$  elements, the sample contains each element seen so far with probability  $S/n$

OCTOBER

SUNDAY



08

281-084 • WK 41

OCTOBER

FRIDAY

3	4	5	6	7	8	9
10	11	12	13	14	15	16
17	18	19	20	21	22	23
24	25	26	27	28	29	30

## Bloom filter

Ex: Create gmail account - Enter user name - everything  
User name already exist)  
↓  
all procedure done quickly  
to check ↓  
↓ { linear search ?  
Binary search ?  
Time consuming process  
↓ use  
Bloom filter

- Bloom filter is a space efficient probabilistic data structure that is used to test whether an element is a member of a set.

for ex checking availability of user name. It is a set membership problem, where the set is a list of all registered user name

- The problem with the bloom filter is that they are probabilistic in nature that means, there might be some false positive results.

False Positive: It might tell that given user name is already taken but actually it's not.

2021

## Working of Bloom filter

- An empty bloom filter is a bit array of m bits. all set to 'zero'. like this -

11 [ 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 ]  
10 0 m

- 12 - we need 'k' number of hash function to calculate the hashes for a given input.

- when we want to add an item in the filter, the bits at k indices  $h_1(x), h_2(x), \dots, h_k(x)$  are set, where indices are calculated using hash functions.

### Example

- Suppose we want to enter 'throw' in the filter, we are using 3 hash functions & a bit array of length 10, all set to 0. initially.

- first we will calculate the hashes as following suppose,

$$h_1('throw') \% 10 = 1$$

$$h_2('throw') \% 10 = 4$$

$$h_3('throw') \% 10 = 7$$

WEDNESDAY

17	18	19	20	21	22	23
24	25	26	27	28	29	30

Now, we will set the bits at indices 1, 4, 7 & 9 to 1.

9

0	1	0	0	1	0	1	0	0
0	1	2	3	4	5	6	7	8

10

11

Again we want to enter ''catch'', similarly we will calculate hashes

12

$$h_1('catch') \cdot 10 = 3$$

$$h_2('catch') \cdot 10 = 5$$

$$h_3('catch') \cdot 1 \cdot 10 = 4$$

13

set the bits at indices 3, 5 & 4 to 1.

14

0	1	0	1	1	1	0	1	0	0
0	1	2	3	4	5	6	7	8	9

Now, if we want to check ''throw'' is present in filter or not. we will do the same process but in reverse order. calculating respective hashes using  $h_1, h_2$  &  $h_3$  & check if all indices set to 1 in the bit array

- If all bits are set then we can say that ''throw'' is 'probably present'.

- If any of the bit at these indices are 0 then ''throw'' is definitely not present'.

2021

TUESDAY

## false positive in Bloom filters:

The question is why we said 'Probably present'  
why this uncertainty.

+ let's take an example:

suppose we want to check whether 'cat' is present or not. we will calculate hashes using  $h_1, h_2 \& h_3$ .

check array

$$h_1('cat') \cdot 10 = 1 \leftarrow \text{set}$$

$$h_2('cat') \cdot 10 = 3 \leftarrow \text{set}$$

$$h_3('cat') \cdot 10 = 7 \leftarrow \text{set}$$

If we check the bit array, bits at these indices are set to 1 but we know that 'cat' was never added to the filter. Bit at index 1,  $\&$  7 was set when we added 'throw' & bit 3 was set when we added 'catch'.

Control probability of getting False Positive by controlling size of Bloom filter.

↓ probability of false positive  $\Rightarrow$  ↑ no. of hash functions.

MONDAY

24 25 26 27 28 29 30

probability of false positive

$$P = \left( 1 - \left[ 1 - \frac{1}{m} \right]^{kn} \right)^k$$

where,  $m$  = size of bit array $n$  = no. of expected elements to be inserted $k$  = no. of hash functions.

$$\text{Size of Bit array} = m = \frac{-n \log P}{(\log 2)^2}$$

$$\text{Optimal no. of hash functions : } k = \frac{m \log 2}{n}$$

Interesting properties of bloom filters

- Bloom filters never generate false negative result i.e. telling you that a username does not exist when it actually exists.
- Deleting elements from filter is not possible because, if we delete a single element by clearing bits at indices generated by  $k$  hash functions, it might cause deletion of few other elements.

1	2	3	10	11	12	13	14
8	9	17	18	19	20	21	
15	16	24	25	26	27	28	
22	23						
29	30						

TUESDAY

## Counting Distinct Elements in a Stream

### Flajolet - Martin Algorithm (FM Algorithm)

- FM algorithm approximates the no. of unique objects in a stream or a database in one pass.
- If the stream contains n elements with m of them unique ; this algorithm runs in  $O(n)$  time & needs  $O(\log m)$  memory.
- The FM algorithm is an algorithm for approximating the no. of distinct elements in a stream with a single pass.
- Space consumption logarithmic in the maximal no. of possible distinct elements in the stream.

#### Example

- Determine the distinct elements in the stream using FM algorithm.

Input stream of integers  $x = 1, 3, 2, 1, 2, 3, 4, 3, 1, 2, 3, 1$

Hash function,  $h(x) = 6x + 1 \bmod 5$

2021

MONDAY

17	18	19	20	21	22	23
24	25	26	27	28	29	30

Step-1 : calculating Hash function [ h(x) ]

$$h(x) = 6x + 1 \bmod 5$$

$$\begin{aligned} h(1) &= 6(1) + 1 \bmod 5 \\ &= 7 \bmod 5 \end{aligned}$$

$$h(1) = 2$$

similarly, calculating hash function

for the remaining  
input stream.

$$h(1) = 2 \quad h(4) = 0$$

$$h(3) = 4 \quad h(3) = 4$$

$$h(2) = 3 \quad h(1) = 2$$

$$h(1) = 2 \quad h(2) = 3$$

$$h(2) = 3 \quad h(3) = 3$$

$$h(3) = 4 \quad h(1) = 2$$

3

Step-2 : Write Binary Equivalent for hash function

$$h(1) = 2 = 010$$

$$h(3) = 4 = 100$$

$$h(2) = 3 = 011$$

$$h(1) = 2 = 010$$

$$h(2) = 3 = 011$$

$$h(3) = 4 = 100$$

$$h(4) = 0 = 000$$

$$h(3) = 4 = 100$$

$$h(1) = 2 = 010$$

$$h(2) = 3 = 011$$

$$h(3) = 3 = 011$$

$$h(1) = 2 = 010$$

Step-3 : Finding Trailing zeros

Now, write the count of trailing zeros in each hash function bit

2021

1	2	3	4	5	6	7
8	9	10	11	12	13	14
15	16	17	18	19	20	21
22	23	24	25	26	27	28
29	30					

OCTOBER

Not

trailing zeros - if  
bit  
change  
then  
trailing  
zero  
consider

SATURDAY

$$h(1) = 2 = 010 = 1$$

$$h(4) = 0 = 000 = 0$$

$$h(3) = 4 = \underline{100} = 2$$

$$h(3) = 4 = 100 = 2$$

$$h(2) = 3 = 011 = 0$$

$$h(1) = 2 = 010 = 1$$

$$h(1) = 2 = 010 = 1$$

$$h(2) = 3 = 011 = 0$$

$$h(2) = 3 = 011 = 0$$

$$h(3) = 4 = 100 = 2$$

$$h(3) = 4 = 100 = 2$$

$$h(1) = 2 = 010 = 1$$

Step-4

write the value of maximum no. of trailing zeros.

- The value of  $\sigma = 2$

- The distinct value  $R = 2^{\sigma}$   
 $= 2^2$   
 $\boxed{= 4}$

Hence, there are 4 distinct elements 1, 3, 2, 4

10

SUNDAY

# PDF Created Using



## Camera Scanner

Easily Scan documents & Generate PDF



<https://play.google.com/store/apps/details?id=photo.pdf.maker>