

Bloom filters

Ch. 8.4

Bloom filters

False positive (filter error) may occur.

Rejection is always accurate.

- When?
 - Returning “Maybe” and “No” as answers is acceptable.

- What?

- Bit array
- | | | | | | |
|---|---|---|---|-----|-------|
| 0 | 1 | 2 | 3 | ... | $m-1$ |
| 0 | 1 | 0 | 0 | ... | 1 |

- Uniform and independent hash functions f_1, f_2, \dots, f_h

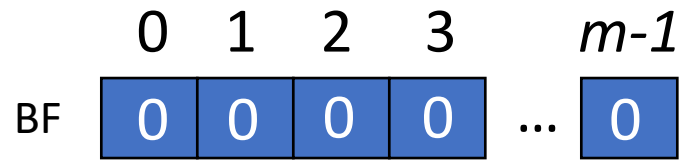
$0 \leq f_i(k) \leq m-1$, where k is key and $i = 1, 2, \dots, h$

- Operations:

- Insert an element into the set
- Member: Is the element in the set?

Operation: Insertion

Given m bits of memory BF and h hash functions.

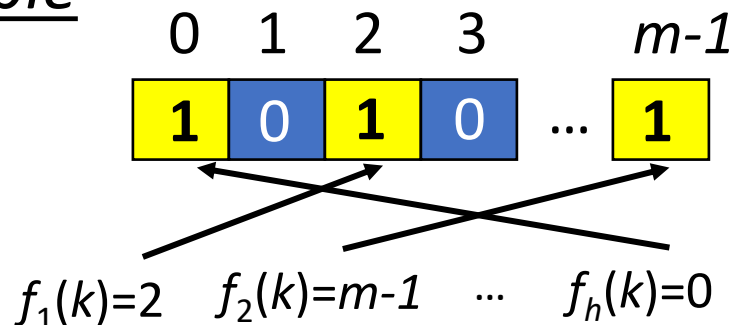


f_1, f_2, \dots, f_h

$$0 \leq f_i(k) \leq m-1$$

- Initialize all m bits to be 0.
- To insert key k , set bits $f_1(k), f_2(k), \dots, f_h(k)$ to be 1.

Example



Signature of key k

Example

Given $m = 11$ (Normally, very larger m is used.)

$h = 2$ (2 hash functions)

- $f_1(k) = k \bmod m$
- $f_2(k) = (2k) \bmod m$

- Operation 1: Insert $k = 15$

$$f_1(15) = 15 \bmod 11 = 4$$

$$f_2(15) = (2 \cdot 15) \bmod 11 = 8$$

0	1	2	3	4	5	6	7	8	9	10
0	0	0	0	0	0	0	0	0	0	0

0	0	0	0	1	0	0	0	1	0	0
---	---	---	---	---	---	---	---	---	---	---

- Operation 2: Insert $k = 17$

$$f_1(17) = 17 \bmod 11 = 6$$

$$f_2(17) = (2 \cdot 17) \bmod 11 = 1$$

0	1	0	0	1	0	1	0	1	0	0
---	---	---	---	---	---	---	---	---	---	---

Operation: Member(k , BF)

Search for key k

- **Any** $\text{BF}[f_i(k)] = 0 \rightarrow k$ is not in the set.
- **All** $\text{BF}[f_i(k)] = 1 \rightarrow k$ may be in the set.

Following the previous example:

$f_1(k) = k \bmod m, f_2(k) = (2k) \bmod m$, where $m=11$

0	1	2	3	4	5	6	7	8	9	10
0	1	0	0	1	0	1	0	1	0	0

The 11-bit bloom filter contains the set $S=\{15, 17\}$.

- Operation 3: Member $k = 15$
 $\text{BF}[f_1(15)] = 1$ and $\text{BF}[f_2(15)] = 1$ return "Maybe"
- Operation 4: Member $k = 6$
 $\text{BF}[f_1(6)] = 1$ and $\text{BF}[f_2(6)] = 1$ return "Maybe"
Filter error or False positive

Exercise

Given $m = 13$ (size of bit array for the bloom filter BF)

$h = 3$ (number of hash functions)

- $f_1(k) = (3k) \bmod m$
- $f_2(k) = (2k) \bmod m$
- $f_3(k) = k^2 \bmod m$

Q3: Please write out the bit array after inserting 11.

Q4: (Continue of Q3) Please write out the bit array after inserting 1.

Q5: (Continue of Q4) What are the results of $\text{Member}(3, \text{BF})$?

Please reply your answers of Q3-Q5 via the following link:



<https://forms.gle/rzSHma4tLUURKEzw5>

Group members: 2~4 people

Design of bloom filters

- Choose m (filter size in bits)
 - Large m to reduce filter error
- Pick h (number of hash functions)
 - h is too small: Probability of different keys having same signature is high.
 - h is too large: The bloom filter fills with ones quickly.
- Select h hash functions
 - Hash functions should be relatively independent.

Performance analysis

- Assume that a bloom filter with
 - m bits of memory
 - h uniform hash functions
 - u elements
- Consider the i -th bit of the bloom filter
 - Probability to be selected by the j -th hash function $f_j(k)$:
$$P[f_j(k) = i] = 1/m$$
 - Probability of unselected by the j -th hash function $f_j(k)$:
$$P[f_j(k) \neq i] = 1 - 1/m$$
 - Probability of unselected by any of h hash functions:
$$P[f_j(k) \neq i \text{ for } j = 1, \dots, h] = (1 - 1/m)^h$$
 - After inserting u elements, probability of unselected by any of h hash functions:
$$p = (1 - 1/m)^{hu}$$

Performance analysis

m bits of memory

h uniform hash functions

u elements

- Consider the i -th bit of the bloom filter
 - After inserting u elements, probability that bit i remains 0: $p = (1 - 1/m)^{hu}$
 - After inserting u elements, probability that bit i is 1: $1 - p$
- Probability of false positives:
 - Take a random element k and check Member(k , BF).
 - What is the probability that it returns *true*?
 - Answer: The probability that all h bits $f_1(k), \dots, f_h(k)$ in BF are 1
 $f = (1 - p)^h$

Selection of h

False positive rate:

$$f = (1 - p)^h, \text{ where } p = (1 - 1/m)^{hu}$$

Two competing forces:

- Large h :
 - Test more bits for $\text{Member}(k, \text{BF}) \rightarrow$ Lower false positive rate
 - More bits in BF are 1 \rightarrow Higher false positive rate
- Small h :
 - Test fewer bits for $\text{Member}(k, \text{BF}) \rightarrow$ Higher false positive rate
 - More bits in BF are 0 \rightarrow Lower false positive rate

Minimizing false positive rate (1)

- Assume that the filter size m and the number of elements in the filter u are fixed,

h minimizes false positive rate f if $h = (m \ln 2) / u$

Proof:

$$\begin{aligned}\min f &= \min (1 - p)^h \\ &= \min e^{h \ln (1-p)} \\ &= \min h \ln (1 - e^{-hu/m}) \\ &= \min h \ln (1 - a^{-h})\end{aligned}$$

$$x = e^{\ln x}$$

Using approximation

$$p = (1 - 1/m)^{hu} \approx e^{-hu/m}$$

$$\text{Let } a = e^{u/m}$$

➡ Differentiating $h \ln (1 - a^{-h})$ with respect to h and setting the result to zeros $\frac{d}{dh} h \ln (1 - a^{-h}) = 0$

$$\begin{aligned}\ln(1 - a^{-h}) + h \frac{a^{-h} \ln a}{1 - a^{-h}} &= 0 \quad \Rightarrow \quad e^{-hu/m} = \frac{1}{2} \quad \Rightarrow \quad h = \frac{m}{u} \ln(2)\end{aligned}$$

Minimizing false positive rate (2)

- Assume that the filter size m and the number of elements in the filter u are fixed,

h minimizes false positive rate f if $h = (m \ln 2) / u$

Proof:

$$\begin{aligned}\min f &= \min (1 - p)^h \\ &= \min e^{h \ln(1-p)} \\ &= \min h \ln(1-p) \\ &\approx \min -\frac{m}{u} \ln(p) \ln(1-p)\end{aligned}$$

$$x = e^{\ln x}$$

Using approximation

$$p = (1 - 1/m)^{hu} \approx e^{-hu/m}$$

$$p = e^{-hu/m} \Rightarrow h = -\frac{m}{u} \ln(p)$$

➡ When $p = 1/2$, the value of f is minimum.

$$\text{➡ } p = e^{-hu/m} = 1/2 \Rightarrow h = -\frac{m}{u} \ln(p) = -\frac{m}{u} \ln\left(\frac{1}{2}\right) = \frac{m}{u} \ln(2)$$

Design of bloom filter

Given m bits of memory and u elements

Choose $h = (m \ln 2) / u$ hash functions

- Probability that some bit i is 1

$$p \approx e^{-hu/m} = 1/2$$

- Expected distribution

$$m/2 \text{ bits } 1, m/2 \text{ bits } 0$$

- Probability of false positives

$$f = (1 - p)^h \approx \left(\frac{1}{2}\right)^h = \left(\frac{1}{2}\right)^{(\ln 2)m/u} \approx 0.6185^{m/u}$$

Limitations of bloom filters

- The naive implementation of the Bloom filter doesn't support the delete operation.
- The false-positives rate can be reduced but can't be reduced to zero.

Exercise

- Given a bloom filter with m bits of memory size and storing u elements. We set $m = 8u$.
 - Q6: Please compute the optimum number of hash functions that minimizes the false positive probability f .
- Q7: (Continue of Q6) Please compute the false positive probability f .

Hint: $f = (1 - p)^h$ $h = \frac{m}{u} \ln(2)$

Please reply your answers of Q6-Q7 via the following link:



<https://forms.gle/rzSHma4tLUURKEzw5>

Group members: 2~4 people

Summary

- What is bloom filter?
- Operations: Insert, Member
- Performance analysis