

Bloom Filter

Probabilistic data structure used by the crawlers to check if a URL has been already crawled or not.

Instead of hashing the URLs and check for collisions, which is too costly in space,

the idea of the Bloom filter is to

create a binary array $B[1, m]$

and map each URL to k

positions in B , using a family

of k hash functions

$(h_i: U \rightarrow \{0, \dots, B-1\} \forall i=1, \dots, k)$.

If $B[h(x)]$ is set to zero, set it to one, otherwise do nothing.

Given a URL v

if $B[h_i(v)] = 1 \forall i=1, \dots, k$

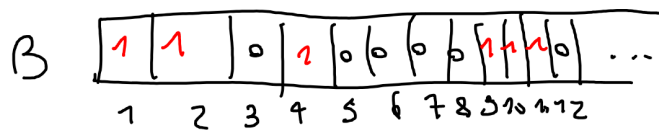
it means that we have already crawled v

As you can see, we don't need to store the whole URL, but only a sort of fingerprint.

However, we might have false positive

EXAMPLE FALSE POSITIVE

Given 3 hash function, a URL that we have yet to crawl (AAA) and a Bloom Filter already populated by other crawled URLs



if $h_1(\text{AAA}) \cap h_2(\text{AAA}) \cap h_3(\text{AAA})$ happens to be equal to a position set to 1, we will wrongly think that AAA was already crawled

What's the probability of false positive?

Assuming that the hashes are distributed perfectly at random (Simple Uniform Hashing)

$$P(B[h_i(x)] = 0) = \left(1 - \frac{1}{m}\right)^{K \cdot n} \approx e^{-\frac{K \cdot n}{m}}$$

where $i = 1, \dots, K$ and x is taken from a set of size n

The probability of a false positive is

$$\begin{aligned} P(\forall i : B[h_i(x)] = 1) &= \\ &= P(B[h_i(x)] = 1)^K = \\ &= \left(1 - P(B[h_i(x)] = 0)\right)^K \approx \\ &\approx \left(1 - e^{-\frac{K \cdot n}{m}}\right)^K \end{aligned}$$

THIS PROBABILITY IS BOUNDED BY

K : LARGE K , LESS FALSE POSITIVE

WHAT'S THE OPTIMAL VALUE FOR K ?

By minimizing the previous probability:

$$\tilde{K} = \ln 2 \cdot \left(\frac{m}{n} \right)$$

Therefore the Bloom filter is good
when:

$$\left(\frac{m}{n} \right) \ll \text{Key length in bits} + \log n$$

SPECTRAL Bloom FILTER

Variation of the Bloom Filter, that was an array of integers C instead of a binary one: each position counts the number of occurrences of an inserted item.

The space occupancy is larger and the error probability is the same, but it also permits deletion and lookup.

Insertion(x)

for h in H :
 $C[h(x)]++$

Deletion(x)

for h in H :
 $C[h(x)]--$

Lookup(x)

return $\min \{ C[h_1(x)], \dots, C[h_k(x)] \}$
(Minimum Selection)