

HW1 - Data Management Algorithms for Decisions Making

Shaza Abdallah - 323032474

November 25, 2025

1 Problem 1: Bloom Filter

1.1

We know from the lecture that the false positive probability of a Bloom filter is

$$\Pr[\text{FP}] = \left(1 - \left(1 - \frac{1}{m}\right)^{kn}\right)^k.$$

we require:

$$\left(1 - \left(1 - \frac{1}{m}\right)^{kn}\right)^k \leq 0.1.$$

after algebraic simplification we get:

$$m \geq \frac{1}{1 - \left(1 - 0.1^{1/k}\right)^{\frac{1}{kn}}}.$$

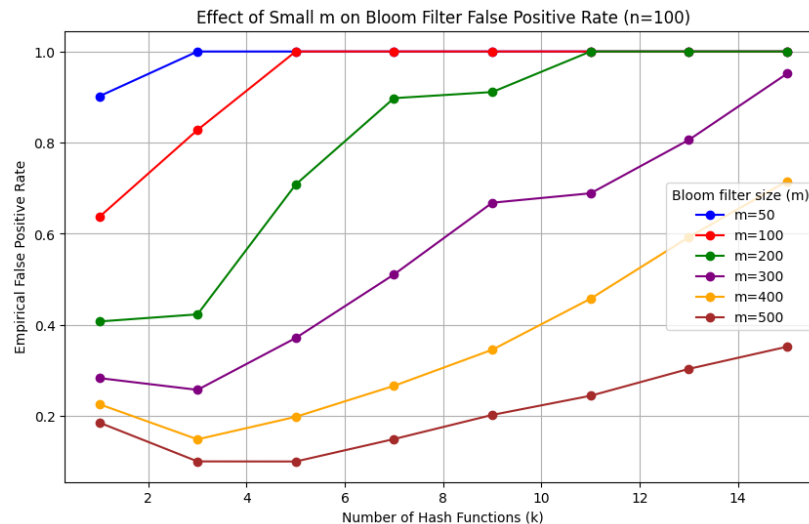
Substituting $n = 1000$ and $k = 5$ gives approximately

$$m \approx 5016.3$$

m is the number of bits so it is an integer, therefore the minimum number of bits is:

$$m = 5017 \text{ bits.}$$

1.2



From the plot, we observe that when the memory m is small, increasing the number of hash functions k is not recommended as the false-positive rate increases. Using large k causes the Bloom filter to become saturated quickly, which significantly increases the false positive rate. This is visible in the plots: the false positive rate rises for high k .

However, when the memory is larger (for example $m \geq 200$), using too few hash functions is also sub-optimal. In particular, for these larger values of m , choosing $k < 3$ increases the false positive rate. In the simulations, $k = 3$ consistently achieves the lowest false positive rate for $m \geq 200$.

This behavior matches the formula

$$\Pr[\text{FP}] = \left(1 - \left(1 - \frac{1}{m}\right)^{kn}\right)^k.$$

As k increases, two opposite effects occur:

- The inner term $(1 - 1/m)^{kn}$ decreases, since we apply more hash functions and thus the entire expression that is raised to the power k increases.
- The entire expression is raised to the power k , decreases as k increases

Because these two effects pull in opposite directions, the false positive probability is not monotonic in k . increasing k is harmful when the memory is limited. Therefore, when m is small, we should decrease k , not increase it but also avoid making k too small, since very small values of k also increase the false positive rate.

1.3

1.3.1 (a)

No, B is not necessarily the same Bloom filter as the one created by inserting all elements of $A_1 \cup A_2$.

Consider the example:

$$A_1 = \{1\}, \quad A_2 = \{2\},$$

using a 5-bit Bloom filter and the hash function

$$h(x) = x \bmod 5.$$

Then the Bloom filters are:

$$B_1 = [0, 1, 0, 0, 0],$$

$$B_2 = [0, 0, 1, 0, 0].$$

Their bitwise AND is:

$$B = B_1 \wedge B_2 = [0, 0, 0, 0, 0].$$

But the Bloom filter created by inserting all elements of

$$A_1 \cup A_2 = \{1, 2\}$$

would be:

$$B_{A_1 \cup A_2} = [0, 1, 1, 0, 0].$$

Clearly,

$$B \neq B_{A_1 \cup A_2}.$$

1.3.2 (b)

No, B does not accurately represent the set $A_1 \setminus A_2$.

In the same example:

$$A_1 \setminus A_2 = \{1\},$$

but the bitwise AND Bloom filter was

$$B = [0, 0, 0, 0, 0],$$

so querying for element 1 would incorrectly return “not present.”

Thus, the Bloom filter obtained by ANDing B_1 and B_2 does not correctly represent the set difference $A_1 \setminus A_2$.

2 Problem 2: Counting Unique Items

2.1

2.1.1 (a)

If we insert n elements into a Bloom filter of length m using k hash functions, then the probability that a specific bit i remains 0 is:

$$\Pr[\text{bit } i = 0] = \Pr[h_k(r) \neq i \ \forall r \in L, \forall k] = \left(1 - \frac{1}{m}\right)^{kn} \approx e^{-kn/m}.$$

Therefore, the expected number of zero bits in the Bloom filter is:

$$c = m \cdot e^{-kn/m}.$$

If we observe c zeros, then solving for n gives:

$$n = \frac{m}{k} \log\left(\frac{m}{c}\right).$$

The average error is 5.7

2.1.2 (b)

(i) Similarly to the above analysis, the probability that bit i is 1 is:

$$\Pr[\text{bit } i = 1] = 1 - \Pr[\text{bit } i = 0] \approx 1 - e^{-kn/m}.$$

Thus, the expected number of one-bits in the Bloom filter is:

$$c = m \left(1 - e^{-kn/m}\right).$$

Solving for n , we obtain:

$$n = -\frac{m}{k} \log\left(1 - \frac{c}{m}\right).$$

(ii) **The average error is 6.35**

(iii) Theoretically, the two estimators are equivalent, since they are derived from the same expression for the Bloom filter's bit statistics. In practice, our simulations show that their average errors are very similar (about 5–6), which is consistent with the theoretical equivalence. The small difference between them likely comes from randomization and from floating-point precision effects when evaluating logarithms on values very close to 0 or 1. So maybe those minor numerical variations produce slightly different estimates in practice.

2.2

2.2.1 (a)

We did a kind of cheating in this implementation, since the number of bits of the hash function should satisfy $l = O(\log n)$, where n is the number of unique items. Thus, we should have an a priori idea of how many unique items could appear in the list. However, since this is a synthetic toy example and the items are random, we believe that actually counting the number of unique items in order to determine the appropriate value of l is a valid approach.

The average error is 9.02

2.2.2 (b)

The analysis is essentially identical: the probability that the r least significant bits are 1 is $1/2^r$, and in expectation we need to observe about 2^r distinct elements before seeing r consecutive ones. Thus, the theoretical behavior is the same as in the trailing-zeros estimator.

The average error is 9.03

2.2.3 (c)

Both estimators are theoretically equivalent, and in our experiments they show very similar average error (approximately 9). However, compared to the previous estimator, they perform worse. This is likely due to their higher variance: the estimator outputs only values of the form 2^r , which makes the estimate less precise.

LLM Use

We did use an LLM in this homework, mainly for translating text into L^AT_EX. Here is an example (and even this section was revised by the LLM :)).

our input: "We did a kind of cheating in this implementation since the number of bits of the hash function should satisfy $l = O(\log n)$ where n is the number of unique items. So we should have an idea about how many unique items could be in the list but since this a synthetic toy example and the items are random we think that actually counting the number of unique items in order to know how much l should be is a validate thing to do. in latex"

LLM output: "We did a kind of cheating in this implementation, since the number of bits of the hash function should satisfy $l = O(\log n)$, where n is the number of unique items. Thus, we should have an a priori idea of how many unique items could appear in the list. However, since this is a synthetic toy example and the items are random, we believe that actually counting the number of unique items in order to determine the appropriate value of l is a valid approach. "

We also used an LLM when working with the trailing estimators. We had an issue with the hash function when truncating its output to l bits, and it helped us correct the implementation (unfortunately we lost the exact prompt).