

We will use the following setting: the universe U is of size M ; the set $S \subset U$ we will do operations on is of size n ; and the table T is of size m . We are considering a Bloom filter with k hash functions $h_i : U \rightarrow T$ for $i \in \{1, 2, \dots, k\}$. Suppose we want to extend Bloom Filters to allow deletions as well as insertions into the underlying set S . We use the following extension referred to as *counting Bloom filters*.

We use an array H where each $H[j]$, $j \in T$, is a b -bit counter instead of a binary bit. Initially each $H(j)$ for $j \in T$ is set to $H[j] = 0$. Each time an element x is inserted into S , $H[h_i(x)]$ is increased by 1 for all $i \in \{1, 2, \dots, k\}$. When we query whether an element x is in S , if $H[h_i(x)] > 0$ for all $i \in \{1, \dots, k\}$ then we report $x \in S$, otherwise we report $x \notin S$. To delete an item x , we first query whether $x \in S$, if we report $x \in S$ then we decrease the counters $H[h_i(x)]$ for all $i \in \{1, \dots, k\}$.

It has been shown that 4 bits per counter is enough for many applications. In this problem we investigate this further. Consider a counting Bloom Filter for a set S of size n , k -hashing functions and m counters. Please answer the following questions:

1. Show that: after n insertions of elements into an empty set S , for each $i \in [0, m-1]$,

$$\Pr\{H(i) \geq j\} \leq 2\left(\frac{enk}{jm}\right)^j.$$

Hint: Consider what is $\Pr\{H(i) = t\}$ first and sum over all $t \geq j$. And use the formula $\binom{n}{i} \leq \left(\frac{ne}{i}\right)^i$.

2. Suppose we choose $k = \ln 2m/n$, argue that after n insertions of elements into an empty set S , the probability that there exists an overflowed counter is tiny in practice. (An upper bound like $c_0 \cdot m$ with $c_0 < 10^{-10}$ will be enough.)
3. Finally, suppose that the size of the counter b is huge, so that no overflow will actually happen. Assume that we first execute n insertions into S . Now, let X_0 be an element in S . After the n insertions, the counters in the hash table are $H[i] = c_i$ for $i \in \{0, 1, \dots, m-1\}$. Let $c_{\min} = \min_{0 \leq i < m} c_i$. Conditioning on the above setting, now we execute t deletions: "DELETE x_j " for $1 \leq j \leq t$, where none of these x_j is in the set S . That is, those t deletions are all invalid. Derive an upper bound of the probability (in terms of c_{\min}) that a false negative occurs when we query "is $X_0 \in S$?" after t deletions.