

Bloom Filters

1. supported operations and use cases

Bloom filters are a data structure that remembers which elements have been seen before and which ones have not, using very little space. For this, it does not store any data directly, but rather consists of many boolean values that are turned to 1 when an element is "inserted". Every element will have a few hash functions that determine which elements of the bloom filter to turn to true, if an element is already true, it will not change how the function works. This also means insert runs at constant time.

To look up an element, one simply has to hash the element and check all the corresponding positions in the bloom filter. This also runs in positive time complexity.

Two potential downsides of bloom filters are the fact that the "delete" operation is not supported, and the existence of "phantom elements" that will cause a false positive rate.¹ Since we are not keeping track of which elements turned a position to True, we might search an element and think it is there because other elements turned all of its hash functions to true.

In summary, the supported operations are:

1. insert
2. search

There are variations of bloom filters that also support deletion, but this exploration will focus on the classic case.

Use cases:

Bloom filters are useful for any situations where memory usage is expensive and the downsides of not being able to delete elements and having the possibility of false positives is not significant. Some examples are Medium keeping track of which articles a reader has already visited so they do not appear in the recommendations, or Google Chrome using Bloom filters to identify malicious URLs (used to).(Roughgarden, 2017)

2. Implementation of a Bloom Filter

```
In [2]: 1 !pip install mmh3
        2 !pip install bitarray
        3
```

Collecting mmh3

Downloading <https://files.pythonhosted.org/packages/fa/7e/3ddcab0a9fcea034212c02eb411433db9330e34d626360b97333368b4052/mmh3-2.5.1.tar.gz> (<https://files.pythonhosted.org/packages/fa/7e/3ddcab0a9fcea034212c02eb411433db9330e34d626360b97333368b4052/mmh3-2.5.1.tar.gz>)

Building wheels for collected packages: mmh3

Running setup.py bdist_wheel for mmh3 ... done

Stored in directory: /root/.cache/pip/wheels/38/b4/ea/6e4e321c625d3320c0c496bf4088371546d8fce5f1dd71b219

Successfully built mmh3

Installing collected packages: mmh3

Successfully installed mmh3-2.5.1

Collecting bitarray

Downloading <https://files.pythonhosted.org/packages/e2/1e/b93636ae36d08d0ee3aec40b08731cc97217c69db9422c0afef6ee32ebd2/bitarray-0.8.3.tar.gz> (<https://files.pythonhosted.org/packages/e2/1e/b93636ae36d08d0ee3aec40b08731cc97217c69db9422c0afef6ee32ebd2/bitarray-0.8.3.tar.gz>)

Building wheels for collected packages: bitarray

Running setup.py bdist_wheel for bitarray ... done

Stored in directory: /root/.cache/pip/wheels/5f/33/3c/1370189fedb8fbb2c3297388cc8cdbe2af0efb48941cf8d526

Successfully built bitarray

Installing collected packages: bitarray

Successfully installed bitarray-0.8.3

```
In [0]: 1 import mmh3 # importing the mmh3 hashing functions
2 from bitarray import bitarray # importing the bitarray structure
3
4 class bloomfilter(object):
5
6     # initializing a bloom filter with size n and n_hash hash functions
7     def __init__(self,length,n_hash):
8         # the size of the bloomfilter is n
9         self.length=length
10        # sets an element equal to n_hash
11        self.n_hash=n_hash
12        # creates a bitarray of size n
13        self.bits=bitarray(length)
14        # sets all the elements in the bitarray to zero
15        self.bits.setall(0)
16
17    # function to insert elements in the bloomfilter
18    def insert(self,elem):
19        # for as many different seeds as there are hash functions
20        for seed in range(self.n_hash):
21            # set the bits at the postion determined by murmurhash equal to 1
22            self.bits[mmh3.hash(elem,seed)%self.length]=1
23
24    # function to search elements in the bloomfilter
25    def search(self,elem):
26        # iterating through n_hash hash function
27        for seed in range(self.n_hash):
28            # using different seeds and adjusting for the length of the list
29            if self.bits[mmh3.hash(elem,seed)%self.length]==0:
30                # return false if the element is not in the bloom filter
31                return False
32            # return true if it is
33            return True
34
35
36    #HC-tag 2
```

```
In [35]: 1 ### Demonstration of bloomfilter ###
2
3 # initialization: a bloomfilter with 2 functions and a bitarray size 10
4 x=bloomfilter(10,2)
5 # inserting a string
6 x.insert("Hello")
7 print("After inserting 'Hello', the Blooom Filter looks like this:",x.bits)
8 # inserting another string
9 x.insert("World")
10 print("After inserting 'World', the bloomfilter looks like this:",x.bits)
11 # searching "Hello" in x
12 print("When searching 'Hello', we get:",x.search("Hello"))
13 # searching "Hello World"
14 print("When searching 'Hello World', we get:",x.search("Hello World"))
```

After inserting 'Hello', the Bloom Filter looks like this: bitarray('1010000000')

After inserting 'World', the bloomfilter looks like this: bitarray('1011000000')

When searching 'Hello', we get: True

When searching 'Hello World', we get: False

Above, we may notice how there is only one additional 1 in the bitarray after we insert the second word. This means one of the hash functions evaluates to the same value. This kind of overlap leads to false positives.

3. Description of the Hash Function used

Murmurhash 3 (mmh3) is a non-cryptographic hashfunction created by Austin Appleby. The name is composed of mu(ltiplication) and r(otation), which are used in the function itself.

I chose this function because it is a fairly standard and efficient algorithm, that passes basic tests like the avalanche test or the chi-squared test.(Appleby, 2011)

To quickly summarize what the function actually does: it takes every 32-bit piece (e.g.4 letters) and shifts them by 8, 16 or 24 bytes depending on their position. We then assign the resulting value to a variable and *multiply* them with different constants, then *rotating* them so the most significant bits become the least significant ones and all the others become more important. There are more details on how to deal with remaining bytes and a final avalanche, that I will not explain here. (Sahib Yar, 2017) It may seem like a complicated algorithm, butthe algorithm compiles down to 52 instructions on a x86 machine (Appleby, 2011), which is incredibly fast for such an efficient algorithm.

```
In [5]: 1 ### Avalanche test ###
2 # shows that with changing a single bit, the output changes dramatically
3 print("mmh3.hash('abd',123)=",mmh3.hash("abd",123))
4 print("mmh3.hash('abe',123)=",mmh3.hash("abe",123))
```

mmh3.hash('abd',123)= 454173339

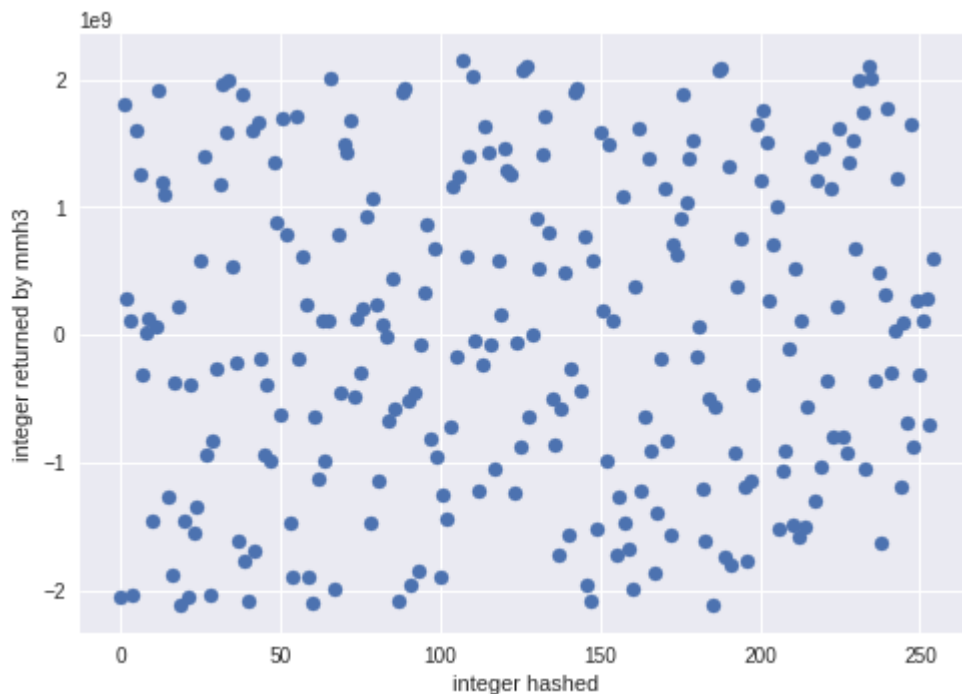
mmh3.hash('abe',123)= -209095228

In [36]:

```

1 ### visual representation of the Chi-squared test ###
2 # the elements are distributed evenly across the range of values (on the y-axis)
3 # this plot also shows that consecutive numbers, represented as bytes,
4 # hash to very different outputs
5 from matplotlib import pyplot as plt
6 import numpy as np
7 y= [mmh3.hash(bytes([i]),2) for i in range(255)]
8 x= [i for i in range(255)]
9 plt.scatter(x,y)
10 plt.xlabel("integer hashed")
11 plt.ylabel("integer returned by mmh3")
12 plt.show()
13
14 # HC-tag 3

```



4. Analysis of scaling

4.1. In terms of memory size as a function of the false positive rate

In my implementation, the false positive rate p is not directly given, but rather determined by the number of elements stored n , the number of hash functions k , and the size of the bitarray m .

$$p = (1 - e^{-\frac{kn}{m}})^k$$

The size of the bitarray is a good measure of the memory size and to relate it to the false positive rate, I can rewrite the function above to return the arraysize m a given false positive rate p , number of inserted elements n and number of hash functions k .

$$m = -kn / \ln(1 - p^{\frac{1}{k}})$$

To optimize further, we can pre-determine the optimal value of k for a given array size and number of elements as

$$k = \frac{m}{n \ln 2}$$

Technically, k needs to be a positive integer, but we will ignore that for now and put it back into the code.

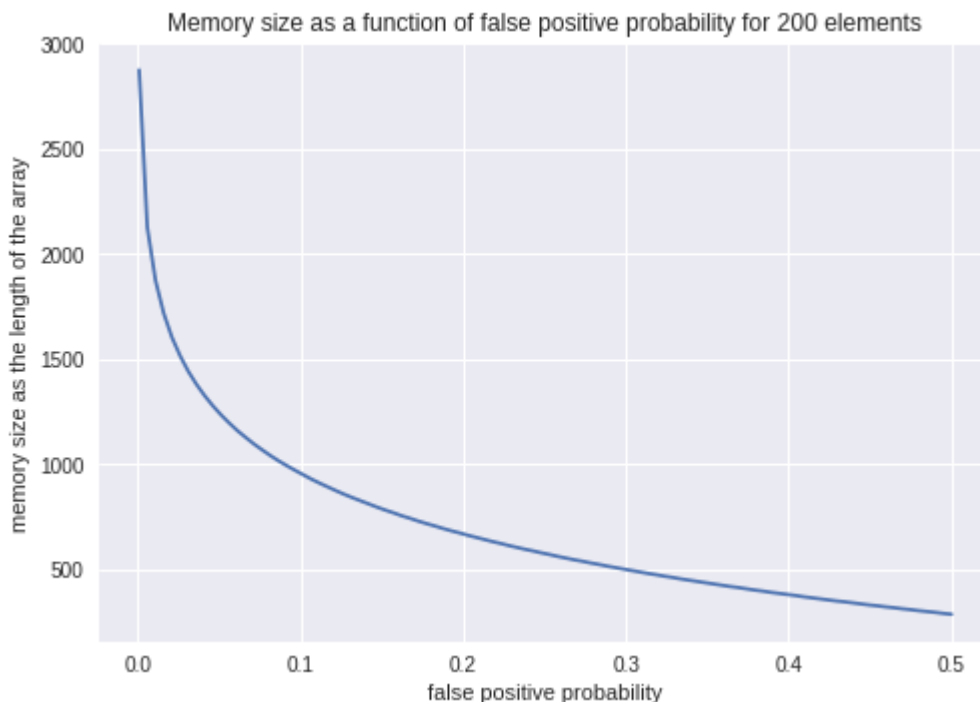
Plugging the function for k back into our function for m , we get that

$$m = -\frac{n \ln p}{(\ln 2)^2}$$

we can now use the function for the size depending on the false positive rate and plot them with a constant number of elements.

To be very thorough, I will additionally define a function that returns the optimal number of hash functions and use my implementation of a bloomfilter to make a version that takes the desired false positive rate and number of elements and creates a corresponding bloomfilter.

```
In [74]: 1 ### plotting memory size as a function of the false positive rate ###
2 import math
3 import numpy as np
4
5 #defining a function that returns the array size given p and n
6 def array_size(p,n):
7     return -(n*math.log(p)/(math.log(2)**2))//1 #//1 to make sure it is an inte
8
9 # setting the input objects to 200
10 n=200
11
12 # plot the array size as a function of p
13 y=[array_size(p,n) for p in np.linspace(0.001,0.5,100)]
14 x=[p for p in np.linspace(0.001,0.5,100)]
15 plt.plot(x,y)
16 plt.title("Memory size as a function of false positive probability for 200 el
17 plt.xlabel("false positive probability")
18 plt.ylabel("memory size as the length of the array")
19 plt.show()
```



4.2 In terms of memory size as a function of the number of items stored

```
In [0]: 1 # define a function that returns the optimal number of hash functions given p
2 def hash_functions(p,n):
3     return 1+(array_size(p,n)/(n*math.log(2)))/1
4
5
6 # define a function that creates a bloomfilter with a specified number of inp
7 # and false positive rate
8 def p_n_bloomfilter(p,n):
9     return bloomfilter(array_size(p,n),hash_functions(p,n))
10
```

4.3 In terms of access time as a function of the false positive rate

The most relevant answer to this question is that a bloomfilter with a predetermined size and number of access functions, the false positive rate will not influence the access time. Instead it will be order $O(k)$, where k is a constant equal to the number of access items.

To make it more interesting, I will however assume that the size of the array changes to adjust to the false positive rate with a set number of elements n .

Using our "hash_functions" function, we can easily plot the access time depending on the false positive rate with a constant value for n .

We notice that the function of k is not continuous since it can only take positive integer values, and decays as p increases.

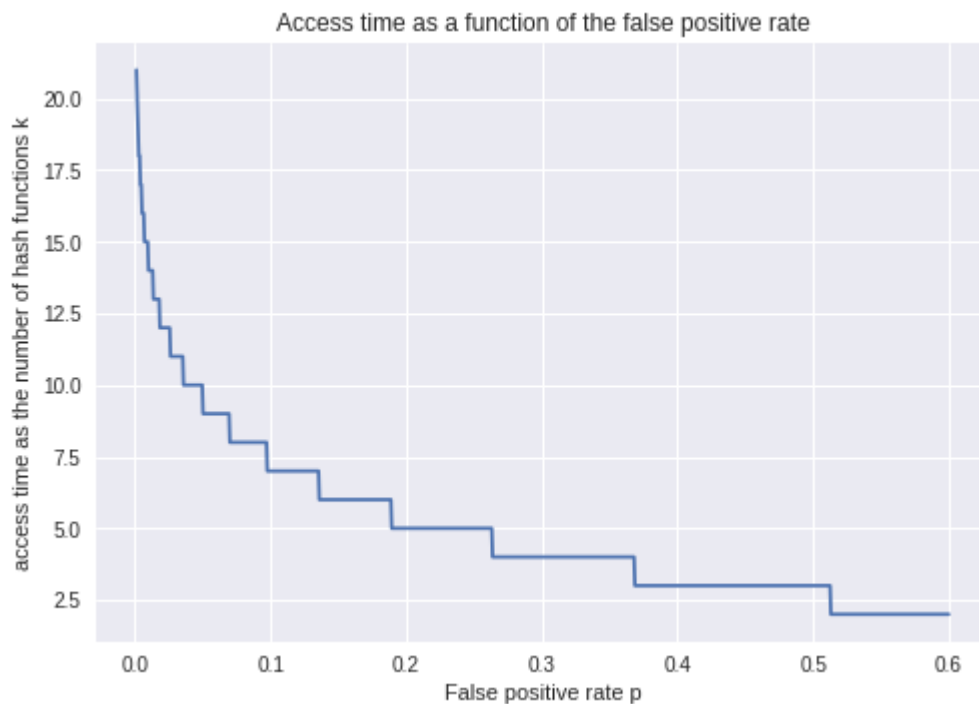
We notice that the value never drops below 1, since we need at least 1 hash-function.

In [85]:

```

1 # set the number of inserted elements equal to 200
2 n=200
3
4 # get the optimal number of hash functions for this number of elements given
5 y=[hash_functions(p,n) for p in np.linspace(0.001,0.6,1000)]
6 # the values of p we iterated over
7 x=[p for p in np.linspace(0.001,0.6,1000)]
8 # make a plot of the two parameters
9 plt.plot(x,y)
10 plt.title("Access time as a function of the false positive rate")
11 plt.xlabel("False positive rate p")
12 plt.ylabel("access time as the number of hash functions k")
13 plt.show()

```



4.4. In terms of access time as a function of the number of items stored

As in 4.3, the complexity of the access operation is still $O(k)$, where k is a constant.

So let's look at the very interesting case of a bloomfilter that automatically adjusts to other requirements:

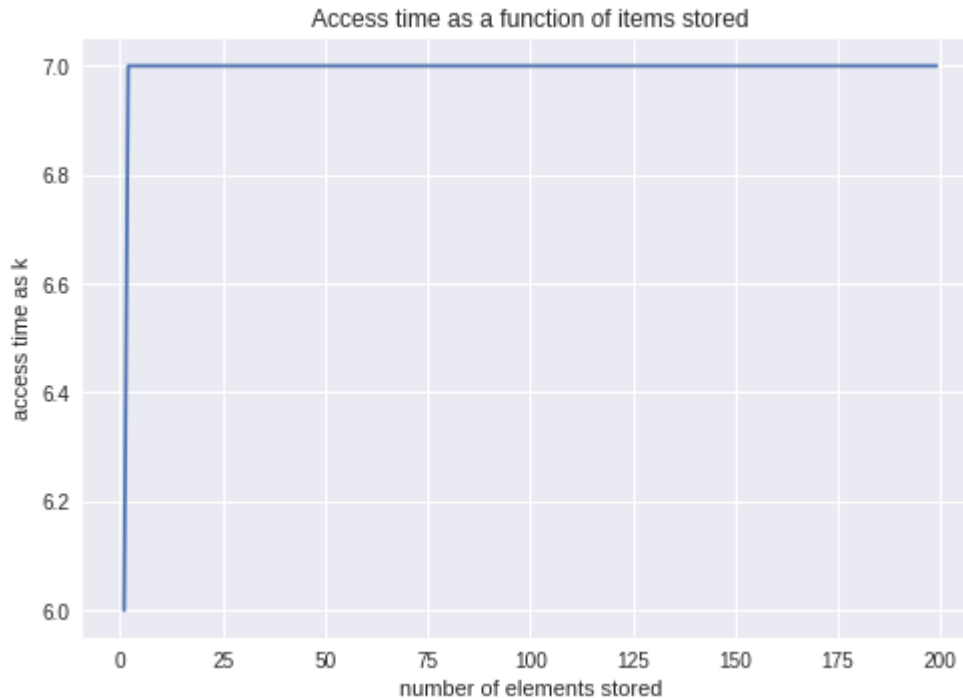
From the collection of functions in part 1, we are tempted to predict the behavior from:

$$k = \frac{m}{n \ln 2}$$

n is in the denominator of the fraction, so we expect a decaying function.

However, To adjust for the size of the array, m will also be adjusted. Under this assumption, we will reach the maximum number of functions for a given value of p rather quickly, after which, the function will be a constant, as expected. The example plotted below is for $p=0.1$


```
In [98]: 1 # set the false positive rate
2 p=0.1
3 # plot the required number of hash functions with array size n
4 y=[hash_functions(p,n) for n in range(1,200)]
5 x=[n for n in range(1,200)]
6 plt.plot(x,y)
7 plt.title("Access time as a function of items stored")
8 plt.xlabel("number of elements stored")
9 plt.ylabel("access time as k")
10 plt.show()
```



5. Comparison of theoretical false positive rate and implementation's false positive rate

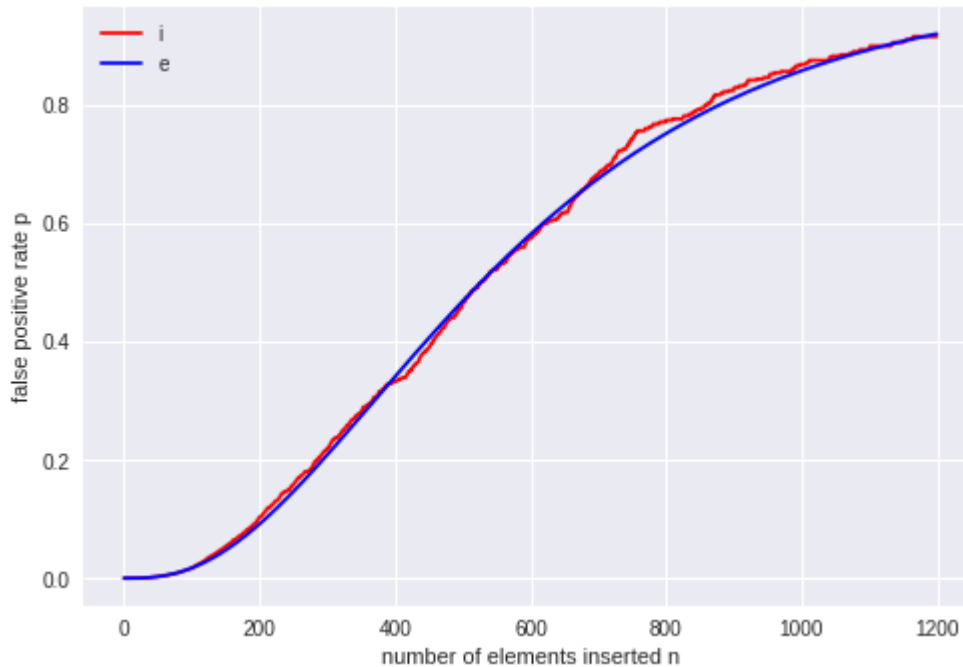
To test the practical false positive rate of the function, we need to generate inputs to test this. First, I will generate a list of 10,000 randomly generated strings containing letters, punctuation and numbers (between 8 and 50 of them). Then I will slowly increase the number of letters that are in the bloom filter and try all the other letters that should not be in the bloom filter. I will keep count of how many of those searches return true and get the percentage of the total. On the same plot, I will show the expected value of the false positive rate with the same array size and number of inserted elements. As we see below, the implementation's false positive rate (red) closely follows the actual false positive rate (blue). Of course, we could also choose m as the x-axis with a constant n .

```
In [0]: 1 import string
2 from random import *
3
4
5 characters = string.ascii_letters + string.punctuation + string.digits
6 words = ["".join(choice(characters) for x in range(randint(8,50))) for i in range(10000)]
```

$$p = (1 - e^{-\frac{kn}{m}})^k$$

```
In [0]: 1 # setting length and number of bloomfilters
        2 m=1000
        3 k=3
        4 # initializing bloomfilter
        5 x=bloomfilter(m,k)
        6 # empty list to collect the false positive rates
        7 y=[]
        8 # iterate through 1200 insertions
        9 for i in range(1200):
        10     x.insert(words[i])
        11     true_count=0
        12     total=1
        13 # measure the false positive rate for each insertion
        14     for j in range(i,len(words)):
        15         if x.search(words[j])==False:
        16             total+=1
        17         else:
        18             true_count+=1
        19             total+=1
        20     y.append(true_count/total)
        21
        22 # calculate the theoretical values over the same range
        23
        24 t=[(1-math.exp(-k*n/m))**k for n in range(1200)]
        25
        26
```

```
In [109]: 1 # plot the implementation's false positive rate as a function of elements
2 # inserted in a bloomfilter of length 1000, with 3 hash functions
3 n=[n for n in range(1200)]
4 plt.plot(n,y,"r")
5 plt.plot(n,t,"b")
6 plt.legend("ie")
7 plt.xlabel("number of elements inserted n")
8 plt.ylabel("false positive rate p")
9 plt.show()
```



Sources:

Appleby, A. (2011). murmurhash. Retrieved November 13, 2018, from <https://sites.google.com/site/murmurhash/> (<https://sites.google.com/site/murmurhash/>)

Roughgarden, T. (2017, January 27). 16 1 Bloom Filters The Basics 16 min - YouTube. Retrieved November 8, 2018, from <https://www.youtube.com/watch?v=zYlxP7F3Z3c> (<https://www.youtube.com/watch?v=zYlxP7F3Z3c>)

Sahib Yar. (2017). Murmur Hash - Explained. Retrieved from <https://www.youtube.com/watch?v=b8HzEZt0RCQ> (<https://www.youtube.com/watch?v=b8HzEZt0RCQ>)

HC-Tags:

1. #algorithms: I implemented an efficient bloomfilter class, clearly annotating my code
2. #professionalism: I went the extra mile in visualizing the hash function as well as all the relations in part 4.
3. #constraints: When introducing bloom filters, I give an explanation of their advantages and disadvantages and why working with the constraints is still beneficial.