

# Com S 435/535 Programming Assignment 1

## 600 Points

Due: Oct 2, 11:59PM

Late Submission Due: Oct 3, 11:59PM(25% Penalty)

In this programming assignment, you will design a Bloom Filter, use it in an application, and will implement the count min sketch data structure.

Description of a programming assignment is not a linear narrative and may require multiple readings before things start to click. You are encouraged to consult instructor/Teaching Assistant for any questions/clarifications regarding the assignment.

Your programs must be in Java, preferably Java 8.1. **You are allowed to work in groups of size 2.** Please do not forget to read the guidelines before you start implementation.

## 1 Bloom Filter

Recall that Bloom Filter is a probabilistic, memory-efficient data structure to store a set  $S$ . You will design the following classes/programs.

- BloomFilterFNV
- BloomFilterRan
- BloomFilterMurmur
- FalsePositives
- BloomDifferential
- NaiveDifferential
- EmpiricalComparison
- CMS

### 1.1 BloomFilterFNV

Recall that to implement a Bloom filter, one needs to choose appropriate hash function(s). In lecture, we discussed two types of hash functions—*deterministic hash functions* and *random hash functions*. For this class, you must use the deterministic hash function 64-bit FNV. Your class should have following constructors and methods.

`BloomFilterFNV(int setSize, int bitsPerElement)`. Creates a Bloom filter that can store a set  $S$  of cardinality `setSize`. The size of the filter should approximately be `setSize * bitsPerElement`. The number of hash functions should be the optimal choice which is  $\ln 2 \times \text{filterSize} / \text{setSize}$ .

`add(String s)`. Adds the string  $s$  to the filter. Type of this method is void. This method should be case-insensitive. For example, it should not distinguish between “Galaxy” and “galaxy”.

`appears(String s)`. Returns `true` if `s` appears in the filter; otherwise returns `false`. This method must also be case-insensitive.

`filterSize()` Returns the size of the filter (the size of the table).

`dataSize()` Returns the number of elements added to the filter.

`numHashes()` Returns the number of hash function used.

In addition to the above methods, you may include additional public/private methods and constructors.

**Remark.** To implement this class, you need  $k$  hash functions. However, 64-bit FNV is a single hash function, that gets a `String` as input and outputs a single hash value (a 64 bit int). To implement the filter, you need to generate  $k$  hash values. How can you do that? Think about it.

## 1.2 BloomFilterMurmur

This class is exactly same as before expect that you will use `murmur` hash function instead of FNV. There are several variants of `murmur`, you must use the 64-bit function implemented at <http://d3s.mff.cuni.cz/~holub/sw/javamurmurhash/MurmurHash.java>

## 1.3 BloomFilterRan

This class is exactly same as before except you can use your own choice of a random hash function. You should use the following random hash function: Pick smallest prime  $p$  that is at least  $tN$ , randomly pick  $a$  and  $b$  from  $\{0, \dots, p-1\}$ , and the hash function is defined as  $(ax + b) \% p$ . Here  $t$  denotes `bitsPerElement` and  $N$  denotes the `setSize`.

## 1.4 FalsePositives

Theoretically, for non-dynamic filters the false positive probability is  $(0.618)^{\text{bitsPerElement}}$ . However, this is under the assumption that the hash functions are picked uniformly at random. In practice, it is not feasible to pick (and store) hash functions uniformly at random. So in practice false positive probability could be bit higher. How can you empirically evaluate the false probability of the filters that you designed above? Design an experiment to evaluate the false probability rate, and implement it. Write a program named `FalsePositives` to estimate the false positive probability of filters—`BloomFilterFNV`, `BloomFilterMurmur`, `BloomFilterRan`.

## 1.5 Application—Differential Files

Bloom filters can be used in some low-level applications that involve file management. An example is that of *differential files*. Consider a large database of records, where each record is of the form  $\langle \text{key}, \text{value} \rangle$ . *key* and *value* could be strings. Suppose that such a data is stored in a single file. Lets name this file as `database.txt`. This file has one record per line. Often we would like to make changes to records. However, it is very inefficient to make changes to the file `database.txt` every time a record is changed (and is prone to errors). A better approach is to create a differential file. Whenever a particular record is changed, then that record is written to a differential file

(lets call this file `differential.txt`), and `database.txt` is not changed. Once the differential file becomes large enough (and when the load on the server is low), then `database.txt` is updated and all contents of `differential.txt` are removed. Typically, the number of records in differential file are less than 10% of the number of records in `database.txt`.

Now consider the query mechanism with such system. User has a **key** and wants to retrieve the **value** associated with that key. For this the system should first check whether **key** appears in the differential file. If it does not appear in the differential file, then the system will access `database.txt`.

Since the number of records in `differential.txt` are much smaller (10%) compared to number of records in `database.txt`, for most of the queries (90% of the time), the system accesses both the files `differential.txt` and `database.txt`. This makes it somewhat inefficient. A solution to speed up this process is to store `differential.txt` (or just the keys of `differential.txt`) in main memory. However, this consumes more memory as number of records in `differential.txt` can be quite large. A solution is via using Bloom filters. Create a Bloom filter of all **keys** in `differential.txt`. Store the filter in main memory. When a key arrives as query, first check if that key appears in the Bloom Filter. If it appears in the filter, then consult `differential.txt`, otherwise access `database.txt` directly. Now, for most keys that do not appear in `differential.txt`, the system does not consult `differential.txt` (unless there is a false positive). This is more efficient.

Your job is to write a program that simulates this. You are given the following files.

- `database.txt`: This is the data base file. Each line of this file is of the following from:  
word1 word2 word3 word4 year1 n1 m1 year2 n2 m2 .....

For example a record would look like

Archbishop had given him 1720 8 6 1727 10 4 1758 20 6 .....

This means the 4-word phrase `Archbishop had given him` appeared 8 times in 6 books in year 1720, and appeared 10 times in 4 books in year 1727 and so on. Here key is the 4-word phrase.

This file was created by processing data from Google Books Ngram Viewer (<http://storage.googleapis.com/books/ngrams/books/datasetv2.html>). This file has nearly 12 Million records.

- `DiffFile.txt`. This is the differential file. This has nearly 1.2 Million records. This is created by picking random records from `database.txt` and making changes to those records.
- You are also given a file named `grams.txt`. This contains all keys that appear in `database.txt`.

Your task is the following. Write a two programs named `BloomDifferential` and `NaiveDifferential` that does the following. Program `BloomDifferential` has

- a method named `createFilter()` that returns a Bloom Filter corresponding to the records in the differential file `DiffFile.txt`.
- a method named `retrieveRecord(String key)` that gets a key as parameter and returns the record corresponding to the record by consulting the Bloom Filter first.

Program `NaiveDifferential` also has a method named `retrieveRecord(String key)` that returns the record corresponding to the key. It does not use Bloom Filter.

How do you compare the performances of the programs `NaiveDifferential` and `BloomDifferential`? Finally, design a program named `EmpiricalComparison` that empirically compares the performances of both the programs.

## 2 Count Min Sketch

Design a class named `CMS` which is the count min sketch to store a multi-set  $S$ . This class will have following constructors and methods.

`CMS(float epsilon, float delta, ArrayList<String> s)` Builds the data structure  $CMS(\epsilon, \delta)$  to store the multi set  $s$ . Note that the number of rows is of the data structure is  $\log_2 \delta$  and number of columns is  $2/\epsilon$ .

`approximateFrequency(String x)` Return the approximate value of  $f_x$  by consulting  $CMS(\epsilon, \delta)$ .

Given a multi set  $S$  consisting of  $N$  items. Recall the following definitions. A  $q$ -heavy hitter is the set of all elements  $x$  for which  $f_x \geq qN$ . A set  $L$  is called *approximate  $\langle q, r \rangle$  Heavy Hitter*, if the following holds: For every  $x \in L$ ,  $f_x \geq qN$  and every  $x$  with  $f_x < rN$  does not belong to  $L$ .

`approximateHeavyHitter(float q, float r)` Return an array list of strings which is a  $\langle q, r \rangle$  heavy hitter. Assume that  $q \geq r + \epsilon$ .

## 3 Report

Write a brief report that includes the following

- For each class that you created, list specifications of all public and private methods that you have written.
- For the classes `BloomFilterFNV`, `BloomFilterMurmur` explain the process via which you are generating  $k$ -hash values, and the rationale behind your process.
- The random hash function that you used for the class `BloomFilterRan`, again explain how you generated  $k$  hash values.
- The experiment designed to compute false positives and your rationale behind the design of the experiment.
- For all the Bloom filter classes report the false probabilities when `bitsPerElement` are 4, 8 and 10. How do false positives depend on `bitsPerElement`? Which filter has smaller false positives? If there is a considerable difference between the false positives, can you explain the difference? How far away are the false positives from the theoretical predictions?
- Describe the how `EmpiricalComparison` is comparing the performances of `BloomDifferential` and `NaiveDifferential`. Explain the rationale behind your design.

- Create CMS data structure by adding words from the file given file `shakespeare.txt`. Add only words whose length is at least 3 and do not add the words “the” or “The” to the data structure. Take  $\epsilon = 1/100$  and  $\delta = 2^{-20}$ .
- Using the above built CMS data structure, compute a set  $L$  that is  $\langle 0.04, 0.03 \rangle$  approximate heavy hitter. Report the following :
  - Number of 0.04 heavy hitters that are in  $L$ .
  - Number of 0.025 heavy hitters that are in  $L$ .
  - Number of items in  $L$  that are not 0.04 heavy hitters.
  - Total number of strings that are added to the data structure, and the total number of distinct strings that are added to the data structure.
  - An estimate of total memory used to store the CMS data structure.

## 4 Guidelines

An obvious choice is to store the bits of the filter in a Boolean array. However, in Java, Boolean array uses a lot more memory. Each cell of the Boolean array may use up to 4 bytes! This defeats the purpose of creating Bloom Filters. Thus your code must use the class `BitSet`. See Java API for more on this class.

Your code should not use any of the Java’s inbuilt functionalities to create hash tables in the Bloom filter classes and the CMS class. For example, your code should not use classes such as `Hashtable` or `HashMap` in `BloomFilterRan`, `BloomFilterMurmur`, `BloomFilterFNV`, and `CMS`. If you wish you may use the method `hashCode()` from the class `String`. This method converts a `String` to an `int`.

**You are allowed to work in groups of size 2.** You are allowed to discuss with your classmates for this assignment. Definition of *classmates*: Students who are taking Com S 435/535 in Fall18. However, You should write your programs without consulting other groups. In your report you should acknowledge the students with whom you discussed and any online resources that you consulted. This will not affect your grade. Failure to acknowledge is considered *academic dishonesty*, and it will affect your grade.

## 5 What to Submit

Please submit all .java files and the report via Canvas. Your report should be in pdf format. Please do not submit .class files. Please zip all .java files and the report, name the file `PA1YourUserID.zip`. **Only one submission per group please.**

Have Fun!