

## CS358: Applied Algorithms

### Assignment 4: Streaming (due 10/20/2021 10PM EDT)

*Instructor: Sam McCauley*

## Instructions

All submissions are to be done through github. This process is detailed in the handout “Handing In Assignments” on the course website. Answers to the questions below should be submitted by editing this document. All places where you are expected to fill in solution are marked in comments with “FILL IN.”

Please contact me at [srm2@cs.williams.edu](mailto:srm2@cs.williams.edu) if you have any questions or find any problems with the assignment materials.

## Problem Description

In this problem, we will be analyzing a very long novel: “In Search of Lost Time,” by Marcel Proust.<sup>1</sup> This novel contains about a million words (including duplicates), and the text file given for this assignment is about 7MB. Nonetheless, we will be using very small streaming data structures to analyze this file with just a single pass over the data—the first using a handful of kilobytes of space, the second using just 32 bytes.

In this assignment, you will be building two data structures.

First, you will build a Count-Min Sketch data structure. All words in “In Search of Lost Time” will be inserted into the Count-Min Sketch. At the end, your data structure will be queried with some of the most common words in the novel: how many times does this word appear? The testing program will compare your output to the actual count of each word; your data structure should always overestimate the count, but give reasonably similar values.

Second, you will build a HyperLogLog data structure. Again, all words in “In Search of Lost Time” will be inserted into it. At the end, your data structure will be queried to find out approximately how many unique words occurred in the novel. HyperLogLog uses an incredibly small amount of space, so it is likely that your data structure will have some error. However, it should usually be reasonably close to the correct value.

INPUT: `test.out` is given three arguments. The first is a text document in ASCII format.<sup>2</sup> The second is a text document, where each line contains a word from the first text document, followed by a space, followed by the number of times that word appears in the

---

<sup>1</sup>This novel is, I understand, very popular and well-regarded. However, it was chosen for this class mostly because its copyright has expired.

<sup>2</sup>As with last time, this means there are no accented characters.

first document. The final argument is an integer denoting the number of unique words in the original text document.

To run your program on “In Search of Lost Time,” you would use the following input:

```
./test.out proust.txt words.txt 36372
```

The following input may be useful for testing:

```
./test.out proustShort.txt wordsShort.txt 125
```

OUTPUT: This assignment is unique in this course in that a single answer is not usually marked as correct or incorrect.<sup>3</sup> Instead, the testing program will output, for each word in the second text file, the actual number of occurrences of the word in the text compared to the number output by your Count-Min Sketch. Furthermore, the testing program will output the number of unique words predicted by your HyperLogLog data structure, compared to the actual number of unique words.

INTERPRETING THE OUTPUT: For the large output, the CMS should generally answer most word queries within 1000 of the correct value. Almost all CMS answers should be within 1500 of the correct value. The HLL overall estimation of the number of words should almost always be between 25000 and 50000.

You should use several different seeds to check that your answers satisfy these bounds.

FUNCTIONS: This assignment is, broadly, structured much like the last assignment. The functions for both data structures already exist, and you must fill them in. The code in `test.c` will perform the above tests using the functions you provide.

`cms.c` and `cms.h` contain the code for the Count-Min Sketch data structure. `hll.c` and `hll.h` contain the code for the HyperLogLog data structure.

Here is a list of the functions and how they are used:

```
void cms_instantiate(Cms* cms)
```

This function is called before any other calls to the `cms` functions. You can think of it like a constructor. It should set constants and allocate memory. `cms` is a pointer to a struct that I found useful (you can change this to not use a struct if you wish). You do not need to edit this function if you don’t want to; the version in the assignment is the version I used.

```
void cms_insert(char* word, int length, Cms* cms)
```

This function inserts a new word (given by `word`) into the filter. `length` is the length of the word, and `cms` is a pointer to the Count-Min Sketch we want to insert to.

`test.c` will insert each word in the first document into the Count-Min Sketch by calling this function. These inserts will all occur before any call to `filter_lookup`.

---

<sup>3</sup>This choice is due to two concerns. First, we’re using randomness, so some error is to be expected. (One could run the program multiple times and perform statistical tests, but that’s very much contrary to the spirit of these structures.) Second, these structures are fairly inconsistent: for example, a cuckoo filter will almost always have approximately the same false positive rate on a large dataset; a CMS or HLL may not.

```
int cms_lookup(char* word, int length, Cms* cms)
```

This function looks up a word (given by `word`) in the sketch pointed to by `cms`. It returns an estimate of how often `word` (which has length `length`) occurs in the first document.

```
void hll_instantiate(Hll* hll)
```

This function is called before any other calls to the hll functions. You can think of it like a constructor. It should set constants and allocate memory. `hll` is a pointer to a struct that I found useful<sup>4</sup> (you can change this to not use a struct if you wish). You do not need to edit this function if you don't want to; the version in the assignment is the version I used.

```
void hll_insert(char* word, int length, Hll* hll)
```

This function inserts a new word (given by `word`) into the HyperLogLog data structure `hll`. `length` is the length of the word, and `hll` is a pointer to the structure we want to insert to.

`test.c` will insert each word in the first document into the HyperLogLog data structure by calling this function.

```
int hll_estimate(Hll* hll)
```

This function asks for an estimate of how many unique words have been inserted into `hll`.

COUNT-MIN SKETCH PARAMETERS: Your CMS should have 4 rows, each of 300 entries. Each entry should be of 32 bits.<sup>5</sup>

HYPERLOGLOG PARAMETERS: Your HyperLogLog data structure should keep track of 32 counters, each of length 8 bits. For 32 counters, the bias constant is .697. (The bias constants are included in `hll.h`.)

---

<sup>4</sup>This is a bit wasteful, unlike the CMS and cuckoo filter. You could easily have the seed as a global constant and just pass around the table of counters rather than storing this struct.

<sup>5</sup>This is wasteful! Unfortunately, 16 bit integers are JUST small enough to barely overflow on this data. Using 18 or 19 bit entries would almost certainly be ideal.

## Questions

**Code** (50 points). Implement a Count-Min Sketch and a HyperLogLog counter, each as described above. You do not need to describe your implementation—but I left some space below in case there's something that you think would be helpful to state.

*Solution.*

**Problem 1** (10 points). Let  $X$  be a positive random variable (i.e.  $X$  only takes on values that are  $\geq 0$ ). Show that

$$\Pr[X \geq e \cdot \mathbb{E}[X]] \leq 1/e.$$

*Hint:* Use the definition of expectation (and the assumption that  $X$  is always positive). It may also be useful to write  $\Pr[X \geq e \cdot \mathbb{E}[X]] = \sum_{i=e \cdot \mathbb{E}[X]}^{\infty} \Pr[X = i]$ .

*Solution.*

**Problem 2** (20 points). Consider a situation where I have a stream of 1,001,000 elements. 1,000,000 of the elements  $a_1, \dots, a_{1000000}$  only appear once, but one element,  $q$ , appears 1000 times throughout the stream.

For each of the following use cases, give the appropriate parameters<sup>6</sup> for a correct Count-Min Sketch with the smallest possible space. Be sure to answer:

- How many rows should I have? (`numRepetitions`)
- How many entries should I have in each row? (`numEntries`)
- How large should each entry be in bits?<sup>7</sup> (The size of each entry in `table`, in bits)

(a) After inserting all stream elements into my CMS, I query an element  $a_i$ , and I receive an answer  $o_i$ , and I query  $q$ , and receive an answer  $o_q$ .

How should I set the parameters of my Count-Min Sketch so that at least 90% of the time,  $o_i \leq o_q$ ? That is to say, how do I set my parameters so that 90% of the time, the CMS accurately returns that  $q$  was more common than  $a_i$ ?

You do not need to give a precise answer; giving parameters that work is OK even if slightly different parameters are more efficient. But, you should explain how you got your answer, and why it attains the required performance.

*Hint:* One way to do this is to come up with constants  $c_q$  and  $c_i$  such that for your parameters you can show that:

1.  $c_i \leq c_q$ , and
2. 90% of the time,  $o_i \leq c_i$  and  $c_q \leq o_q$ .

<sup>6</sup>Here I am asking about the parameters you should use for your *code*: number table entries, number of rows, etc. Determining  $\epsilon$  and  $\delta$  is likely an important part of your answer, but it is not the final solution.

<sup>7</sup>Please give an exact answer, even if C does not support variables of this size.

*Solution.*

(b) After inserting all stream elements into my CMS, I query all elements  $a_1, \dots, a_{1000000}$  to obtain answers  $o_1, \dots, o_{1000000}$ , as well as querying  $q$  to obtain an answer  $o_q$ .

How should I set the parameters of my Count-Min Sketch so that at least 90% of the time,  $\max_i o_i \leq o_q$ ? That is to say, how do I set my parameters so that 90% of the time, the CMS accurately returns that  $q$  was the most common out of all elements queried?

*Hint:* You probably want to use the union bound on part (b) of this question, combining it with the answer to part (a).

*Solution.*

**Problem 3** (10 points). Let's say I create a HyperLogLog structure  $H_1$  for a stream  $a_1, \dots, a_n$  and a second HyperLogLog structure  $H_2$  for a stream  $b_1, \dots, b_n$ . Assume that all  $a_i$  and  $b_i$  are in the same universe  $U$ , and assume that identical parameters and hash functions are used to create  $H_1$  and  $H_2$ .

Describe how to use  $H_1$  and  $H_2$  to create a new HyperLogLog structure  $H_3$  that can estimate the number of unique items in the concatenation of the two streams:  $a_1, \dots, a_n, b_1, \dots, b_n$ . (You should build  $H_3$  using only  $H_1$  and  $H_2$ , without seeing the stream again.)

*Solution.*

**Problem 4** (10 points). The following short questions ask about different situations in which you can use a Count-Min Sketch. For each, say how many entries you would need in each column (in other words, what value of  $\varepsilon$  you would need asymptotically) to obtain the correct answer in a single row with constant probability. In each case, assume that there are  $N$  total insertions into the CMS.

**For each explanation, please give a brief explanation ( $\approx$  one sentence) as to why—and bear in mind that I am only expecting an asymptotic answer.**

(a) A stream consists of two kinds of items: items of the first kind appear once each, and items of the second kind appear twice each. Can a Count-Min Sketch determine if a given item appeared once or twice? How many entries would you need in each row of the CMS in order for it to provide the correct answer with constant probability?

*Solution.*

(b) In a stream of items, one appears the majority of the time (more than half the items are the same). Can the Count-Min Sketch determine the majority item? How many entries would you need in each row of the CMS in order for it to provide the correct answer with constant probability?

*Solution.*