

Asmt 4: Frequent Items

Turn in a hard copy at the start of class:

Wednesday, March 26

10 points

Overview

In this assignment you will explore finding frequent items in data sets, with emphasis on techniques designed to work at enormous scale.

You will use two data sets for this assignment:

- <http://www.cs.utah.edu/~jeffp/teaching/cs5140/A4/S1.txt>
- <http://www.cs.utah.edu/~jeffp/teaching/cs5140/A4/S2.txt>

These data sets each describe a set of 200 lowercase letters, separated by spaces. The order of the file represents the order of the stream. The data sets are small enough that it may be possible to run the algorithms below *by hand* if you have less programming experience. However, it may also be useful (or faster) to implement the algorithms.

As usual, it is highly recommended that you use LaTeX for this assignment. If you do not, you may lose points if your assignment is difficult to read or hard to follow. Find a sample form in this directory: <http://www.cs.utah.edu/~jeffp/teaching/latex/>

1 Streaming Algorithms

A (4 points): Run the Misra-Gries Algorithm (see **L12.2.2**) with $(k - 1) = 9$ counters on streams S1 and S2. Report the output of the counters at the end of the stream.

In each stream, from just the counters, report how many objects *might* occur more than 20% of the time, and which must occur more than 20% of the time.

B (4 points): Build a Count-Min Sketch (see **L12.2.3**) with $k = 10$ counters using $t = 5$ hash functions. Run it on streams S1 and S2.

For both streams, report the estimated counts for objects **a**, **b**, and **c**. Just from the output of the sketch, which of these objects, with probably $1 - \delta = 31/32$, *might* occur more than 20% of the time?

C (2 points): How would these algorithms need to change (to answer the same questions) if each object of the stream was a “word” seen on Twitter, and the stream contained all tweets?

Name one advantage of Count-Min Sketch over the Misra-Gries algorithm.

2 BONUS

The exact heavy-hitter problem is as follows: return *all* objects that occur more than 10% of the time. It cannot return any false positives or any false negatives. In the streaming setting, this requires $\Omega(\min\{m, n\})$ space if there are n objects that can occur and the stream is of length m .

A: (1 point) A 2-Pass streaming algorithm is one that is able to read all of the data in-order exactly twice, but still only has limited memory. Describe a small space algorithm to solve the exact heavy hitter problem (say for 10% threshold).