

Seoul National University

M1522.001400 Introduction to Data Mining

Spring 2016, Kang

Homework 3: Mining Data Streams (Chapter 4)

Due: Apr 13, 09:30 AM

Reminders

- The points of this homework add up to 100.
- Like all homeworks, this has to be done individually.
- Lead T.A.: Jinhong Jung (montecast9@gmail.com)
- Please type your answers in English. Illegible handwriting may get no points, at the discretion of the graders.
- If you have a question about assignments, please upload your question in eTL.
- If you want to use slipdays or consider late submission with penalties, please note that you are allowed one week to submit your assignment after the due date. That is, after Apr 18, we will NOT receive your submission for this assignment.

Remember that:

- Whenever you are making an assumption, please state it clearly

Question 1

Suppose you use a bloom filter for 8 billion bits and 1 billion keys. Calculate the false-positive rate when we use the following number of hash functions. [25 points]

(a) Three hash functions.

(b) Four hash functions.

(Solution) The false positive probability of a bloom filter is $\left(1 - e^{-\frac{km}{n}}\right)^k$ where k is the number of hash functions, m is the number of keys (darts), and n is the number of bits (targets).

(a) $\left(1 - e^{-\frac{3}{8}}\right)^3 = 0.0306$

(b) $\left(1 - e^{-\frac{4}{8}}\right)^4 = 0.0240$

Question 2

Suppose you are given the following stream:

b, c, a, d, a, c, d, b, a, b, a, c, d

Answer the following questions. [25 points]

(a) What is the surprise number (second moment) of the stream?

(b) What is the third moment of the stream?

(Solution) In a stream, the k-th moment is $\sum_{i \in A} (m_i)^k$ where A is a set of N values, and m_i is the number of times value i occurs in the stream.

(a) Second moment: $4^2 + 3^2 + 3^2 + 3^2 = 43$

(b) Third moment: $4^3 + 3^3 + 3^3 + 3^3 = 145$

Question 3

Suppose we are given the stream of *Question 2*, to which we apply the Alon-Matias-Szegedy (AMS) Algorithm to estimate the surprise number. Suppose we keep four variables X_1, X_2, X_3 and X_4 . Assume that we randomly pick the 2nd position for X_1 , the 5th position for X_2 , the 7th position for X_3 , and the 8th position for X_4 to define the four variables. After looking at the stream, answer the following questions using the AMS algorithm. For simplicity, assume that you know the length of the stream (in this case, it's 13) in advance. [25 points]

- (a) What are X_i . *element* and X_i . *value* for the four variables?
- (b) What is the estimated surprise number?

(Solution)

- (a) Based on the AMS algorithm, the element and the value of each variable are as follows:

variable	X_i . element	X_i . value
X_1	c	3
X_2	a	3
X_3	d	2
X_4	b	2

- (b) The surprise number is the average of $f(X_i)$ where $f(X_i) = n(2 \times X_i$. *value* $- 1)$ (i.e., $S = \frac{1}{k} \sum_{i=1}^k f(X_i)$) The number is 52.

variable	$f(X_i)$
X_1	$13 \times (2 \times 3 - 1) = 65$
X_2	$13 \times (2 \times 3 - 1) = 65$
X_3	$13 \times (2 \times 2 - 1) = 39$
X_4	$13 \times (2 \times 2 - 1) = 39$

Question 4

Suppose the window is as shown in Figure 1. Estimate the number of 1's the last k positions, for $k =$ (a) 5 (b) 15. In each case, how far off the correct value is your estimate? [25 points]

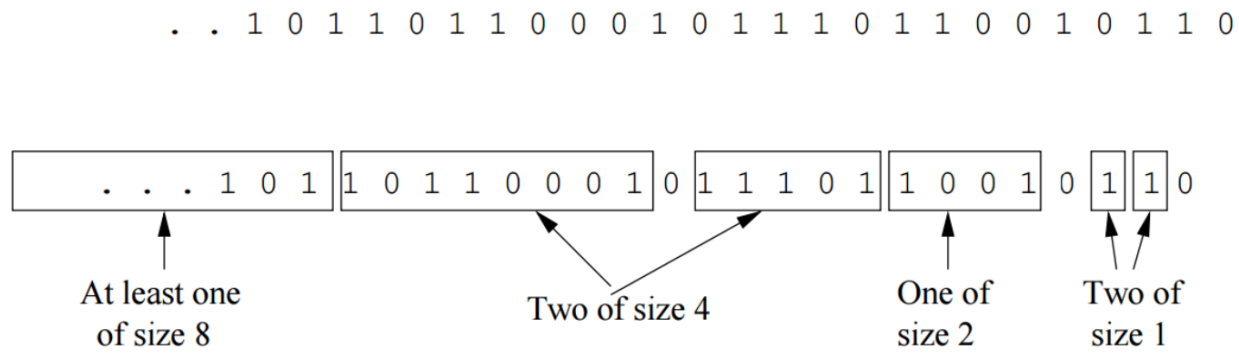


Figure 1. A bit-stream divided into buckets following the DGIM rules

(Solution) In DGIM framework, the number of 1's the last k positions is sum of sizes of more recent buckets except the last block, which includes the k -th bit, plus the half of the size of the last block.

(a) $k = 5$

- a. Correct value: 3
- b. Estimation: $1 + 1 + \left(\frac{2}{2}\right) = 3$
- c. Error: $|3 - 3| = 0$

(b) $k = 15$

- a. Correct value: 9
- b. Estimation: $1 + 1 + 2 + 4 + \left(\frac{4}{2}\right) = 10$
- c. Error: $|10 - 9| = 1$