

Mining Data Streams II

Alexander Schönhuth



Bielefeld University
June 21, 2023

TODAY

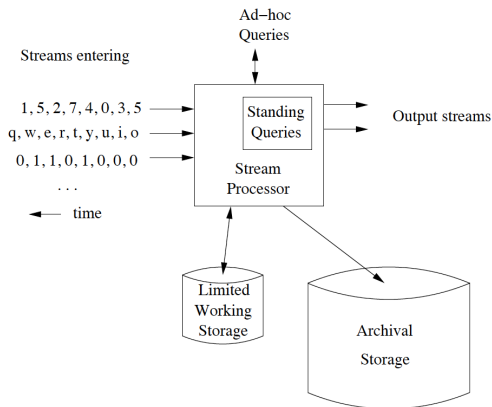
Mining Data Streams II: Overview

- ▶ Counting Ones in a Window:
 - ▣ Datar-Gionis-Indyk-Motwani algorithm
- ▶ Decaying Windows

Learning Goals: Understand these topics and get familiarized

Counting Ones in a Window
The Datar-Gionis-Indyk-Motwani Algorithm

DATA STREAM MANAGEMENT SYSTEM



A data stream management system

Adopted from mmds.org

DATA STREAM QUERIES

Issues

- ▶ Streams deliver elements rapidly: need to act quickly
- ▶ Thus, data to work on should fit in main memory
- ▶ New techniques required:
 - ☞ Compute approximate, not exact answers
 - ☞ Hashing is a useful technique

COUNTING ONES IN WINDOW: PROBLEM

- ▶ Let $x_i \in \{0, 1\}$, earliest = leftmost = x_1 (unlike before)

$$x_1, \dots, x_t, \underbrace{x_{t+1}, \dots, x_{t+N-k+1}, \dots, x_{t+N}}_{\text{window of length } N} \quad (1)$$

- ▶ *Situation:*
 - ▶ We have a window of length N on a binary stream
 - ▶ Query: “how many ones are there in the last $k \leq N$ bits?”
 - ▶ We cannot afford to store entire window
 - ▶ Approximate algorithms required
- ▶ Present solution for binary streams first
- ▶ Thereafter extension for summing numbers

THE COST OF EXACT COUNTS

- ▶ One needs to store N bits to answer count-one-queries for arbitrary $k \leq N$:
 - ▶ Assume one could use less than N bits
 - ▶ We need 2^N different representations to represent all possible 2^N bit strings of length N
 - ▶ Since we use less than N bits, there are two different bit strings $w \neq x$, for which we use the same representation
 - ▶ Let k be the first bit from the right where w and x disagree
 - ▶ *Example:*
 - ▶ For $w = 0101, x = 1010$, we have $k = 1$
 - ▶ For $w = 1001, x = 0101$, we have $k = 3$

THE COST OF EXACT COUNTS

- ▶ One needs to store N bits to answer count-one-queries for arbitrary $k \leq N$:
 - ▶ Let k be the first bit from the right where w and x disagree
 - ▶ *Example:*
 - ▶ For $w = 0101, x = 1010$, we have $k = 1$
 - ▶ For $w = 1001, x = 0101$, we have $k = 3$
 - ▶ So the counts of ones in the window of length k for w and x differ
 - ▶ But because we use identical representations for w and x , we will output the same count
 - ▶ This proves that one needs the full N bits to represent bit strings for exact count-one-queries.

THE DATAR-GIONIS-INDYK-MOTWANI ALGORITHM

► *Situation:*

- We consider a binary stream: elements are *bits*
- Let each element of the stream have a *timestamp*
- The first, *leftmost* element has timestamp 1, the second leftmost has timestamp 2, and so on; i is timestamp for x_i

$$\underbrace{x_1, \dots, x_t, x_{t+1}, \dots, x_M}_{\text{timestamps: } 1, \dots, t, t+1, \dots, M} \quad (2)$$

- *Goal:* We like to count the ones among the N most recent (rightmost) elements/bits

THE DATAR-GIONIS-INDYK-MOTWANI ALGORITHM

- ▶ *Algorithm:* Divide window into *buckets*, contiguous bit substrings
- ▶ *Bucket Representation:* Store
 - ▶ The timestamp (TS) of its right end (figure: $t - 7$), and
 - ▶ The *size* of the bucket, as the number of 1's in the bucket
 - ▶ The size is supposed to be a power of 2 (figure: 2^2)

$$\begin{array}{c} \overbrace{0 [1011 \underbrace{1}_{x_{t-7}}] 011000 \underbrace{1}_{x_t}}^N \\ \underbrace{\hspace{10em}}_{\substack{\text{size } 4=2^2 \\ \text{TS: } t-7}} \end{array} \quad (4)$$

- ▶ *Bucket Space Requirements:*
 - ▶ Storing buckets: (TS, size)
 - ▶ TS between 0 and $N - 1$, so requires $\log_2 N$ bits
 - ▶ Size: Storing $\log_2 j$, $0 \leq j \leq \log_2 N$ amount to $\log_2 \log_2 N$ bits
 - ▶ Requires $O(\log N)$ bits overall

THE DATAR-GIONIS-INDYK-MOTWANI ALGORITHM

Key Ideas / Considerations

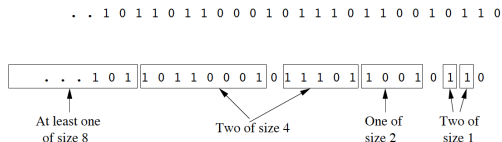
- ▶ Number of buckets representing a window must be small
- ▶ Estimate number of 1's in last k bits by exploiting known (because stored) bucket structure
- ▶ For any k , estimate has error of no more than 50%
- ▶ How to re-establish DGIM Bucket Rules quickly, on new bits arriving?

THE DATAR-GIONIS-INDYK-MOTWANI ALGORITHM

Storage Requirements Overall

- ▶ Each bucket represented using $O(\log N)$ bits (see before)
- ▶ Let 2^j be size of largest bucket: $2^j < N$ implies $j \leq \log_2 N$
- ▶ So there are at most 2 buckets of sizes $2^j, j = \log_2 N, \dots, 1$
- ▶ This means that there are $O(\log N)$ buckets overall
- ▶ $O(\log N)$ buckets of $O(\log N)$ bits: $O(\log^2 N)$ space overall

THE DATAR-GIONIS-INDYK-MOTWANI ALGORITHM



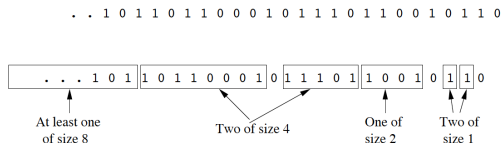
Bit stream divided into buckets following DGIM rules

From mmds.org

Answering Queries

- ▶ Let $1 \leq k \leq N$: how many 1's are among the last k bits?
- ▶ *Answer:*
 - ▶ Find leftmost (= with earliest timestamp) bucket b containing some of last k bits
 - ▶ *Estimate:* Sum of sizes of buckets right of b plus half the size of b

THE DATAR-GIONIS-INDYK-MOTWANI ALGORITHM



Bit stream divided into buckets following DGIM rules

From mmds.org

Example

- ▶ Let $k = 10$: how many 1's are among 0110010110?
- ▶ Let t be timestamp of rightmost bit
- ▶ Buckets of timestamps $t - 1, t - 2$ and size 1 fully included in k rightmost bits
- ▶ Bucket of size 2 with timestamp $t - 4$ is also included
- ▶ Bucket of size 4 with timestamp $t - 8$ is only partially included
- ▶ Estimate: $1 + 1 + 2 + (1/2 \times 4) = 6$, one more than true count

DGIM: ERROR OF ESTIMATE

Case 1: estimate is less than c

- ▶ Let c be true count; let leftmost bucket b be of size 2^j
- ▶ *Worst case*: all 1's in b are among k most recent bits
- ▶ *Worst case example*:

$$\begin{array}{c} \text{k most recent bits} \\ \overbrace{[101101101101] 0 [101\dots]} \\ \underbrace{\hspace{10em}} \\ \text{leftmost bucket } b \end{array} \quad (5)$$

- ▶ Left timestamp of leftmost bucket is $t - k + 1$
 - ▶ All ones of left most bucket (here: 8) belong to true count
 - ▶ Estimate counts only half of them (here: 4)
- ▶ Because $c \geq 2^j$, error is at most half of c :

$$\frac{\text{estimate}}{\text{true count}} = \frac{c - 2^{j-1}}{c} = 1 - \frac{2^{j-1}}{c} \stackrel{c \geq 2^j}{\geq} 1 - \frac{2^{j-1}}{2^j} = \frac{1}{2} \quad (6)$$

DGIM: ERROR OF ESTIMATE

Case 2: estimate is larger than c

- ▶ Let c be true count; let leftmost bucket b be of size 2^j
- ▶ Worst case: only rightmost bit of b is among k most recent bits, and
- ▶ There is only one bucket for each of sizes $2^{j-1}, \dots, 1$
- ▶ Worst case example:

$$\left[10110110110 \overbrace{1}^{t-k+1} \right] 0 \left[101011 \right] 0 \left[1001 \right] 0 \left[1 \right] \overbrace{0}^t \quad (7)$$

k most recent bits only one bucket of each size

- ▶ Timestamp of leftmost bucket is $t - k + 1$
- ▶ Estimate counts half of ones (here: 4); true count is 1
- ▶ That yields $c = 1 + 2^{j-1} + \dots + 1 = 1 + 2^j - 1 = 2^j$
- ▶ Estimate is $2^{j-1} + 2^{j-1} + \dots + 1 = 2^{j-1} + 2^j - 1$, so
- ▶ Error $\frac{2^{j-1} + 2^j - 1}{2^j}$ is no greater than 50% of true count

MAINTAINING DGIM RULES

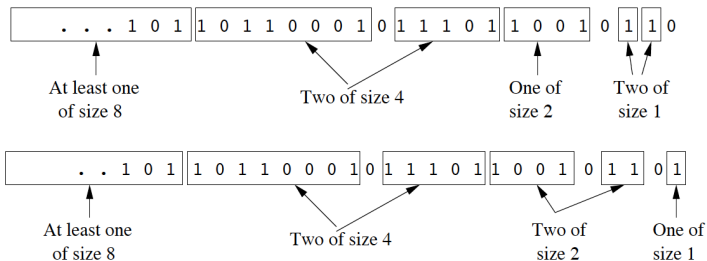
Upon a new bit with timestamp t having arrived:

- ▶ Check timestamp s of leftmost bucket b :
 - ▶ if $s \leq t - N$, drop b from list of buckets
- ▶ If the new bit is 0, do nothing
- ▶ If the new bit is 1, do
 - ▶ Create new bucket with timestamp t and size 1
 - ▶ On increasing size, while there are three buckets of the same size, do
 - ▶ keep the rightmost bucket of that size as is
 - ▶ join the two left buckets into one of double the size
 - ▶ where the timestamp is that of the rightmost bit
 - ▶ *For example:* joining the two left of the three buckets of size 1 into a bucket of size 2 may create a third bucket of size 2, and so on
- ▶ *Runtime:* Need to look at $O(\log N)$ buckets, joining is constant time, so processing new bit requires $O(\log N)$ time overall

THE DATAR-GIONIS-INDYK-MOTWANI ALGORITHM

PART VI

. . 1 0 1 1 0 1 1 0 0 0 1 0 1 1 1 0 1 1 0 0 1 0 1 1 0



Buckets following DGIM rules (top), with new 1 arriving (bottom)

From mmds.org

DGIM ALGORITHM: REDUCING THE ERROR

- ▶ For some $r > 2$, allow either r or $r - 1$ buckets of the same size
- ▶ Allow this for all but size 1 and largest size, whose numbers may be any of $1, \dots, r$
- ▶ Compute estimate as before
- ▶ Extend maintaining the DGIM Bucket Rules in the obvious way
- ▶ *Recall:* largest error $\frac{2^{j-1}+2^j-1}{2^j}$ was made when only one 1 from leftmost bucket b was within window

DGIM ALGORITHM: REDUCING THE ERROR

- ▶ *Recall:* largest error $\frac{2^{j-1}+2^j-1}{2^j}$ was made when only one 1 from leftmost bucket b was within window
- ▶ *New error:*
 - ▶ True count is at most $1 + (r - 1)(2^{j-1} + \dots + 1) = 1 + (r - 1)(2^j - 1)$
 - ▶ Estimate is $2^{j-1} + (r - 1)(2^j - 1)$, difference between estimate and true count is $2^{j-1} - 1$, so fractional error is

$$\frac{2^{j-1} - 1}{1 + (r - 1)(2^j - 1)}$$

which is upper bounded by $1/2(r - 1)$

- ▶ Picking large r can limit error to any $\epsilon > 0$

DGIM ALGORITHM: EXTENSIONS

- ▶ DGIM can be extended to integers instead of bits
- ▶ Question is to estimate the sum of last $k \leq N$ integers from a window of N integers overall
- ▶ However, DGIM cannot be extended to streams containing negative integers
- ▶ Consider case of integers in range of 0 to $2^m - 1$, represented by m bits
- ▶ *Example: $m = 3$, integers 0 to 7*

integer stream :	2	4	3	1	6	7	...	
bit stream :	010	100	011	001	110	111	...	(8)

DGIM ALGORITHM: EXTENSIONS

- ▶ Consider case of integers in range of 0 to $2^m - 1$, represented by m bits
- ▶ *Example:* $m = 3$, integers 0 to 7

$$\begin{array}{rcccccccc} \text{integer stream :} & 2 & 4 & 3 & 1 & 6 & 7 & \dots & \\ \text{bit stream :} & 010 & 100 & 011 & 001 & 110 & 111 & \dots & \end{array} \quad (9)$$

- ▶ *Solution:*
 - ▶ Treat each bit of integers as separate stream
 - ▶ *Example from above:*

$$\begin{array}{rcccccccc} c_0 : & \text{rightmost bit} & 0 & 0 & 1 & 1 & 0 & 1 & \dots \\ c_1 : & \text{middle bit} & 1 & 0 & 1 & 0 & 1 & 1 & \dots \\ c_2 : & \text{leftmost bit} & 0 & 1 & 0 & 0 & 1 & 1 & \dots \end{array} \quad (10)$$

DGIM ALGORITHM: EXTENSIONS

► *Solution:*

- Treat each bit of integers as separate stream
- *Example from above:*

$$\begin{array}{llllllll} c_0 : & \text{rightmost bit} & 0 & 0 & 1 & 1 & 0 & 1 & \dots \\ c_1 : & \text{middle bit} & 1 & 0 & 1 & 0 & 1 & 1 & \dots \\ c_2 : & \text{leftmost bit} & 0 & 1 & 0 & 0 & 1 & 1 & \dots \end{array} \quad (11)$$

- Apply DGIM algorithm to each of m streams \Rightarrow estimate c_i for i -th stream
- *Overall estimate:*

$$\sum_{i=0}^{m-1} c_i 2^i$$

- If error is at most ϵ for all i , overall error is also at most ϵ

Most Common Elements
Decaying Windows

DECAYING WINDOWS: MOTIVATION

- ▶ *Stream*: Movie tickets purchased all over the world
- ▶ *Goal*: Listing currently most “popular” movies
- ▶ *Currently popular*:
 - ▶ Movie that sold plenty of tickets years ago not to be listed
 - ▶ Movie that sold $2n$ tickets last week, for large n , currently popular
 - ▶ Movie that sold n tickets in last 10 weeks is even more popular
 - ▶ How to grasp that idea?

DECAYING WINDOWS: MOTIVATION

- ▶ *Stream*: Movie tickets purchased all over the world
- ▶ *Goal*: Listing currently most “popular” movies
- ▶ *Possible solution*:
 - ▶ One bit stream for each movie
 - ▶ i -th bit in a movie stream is 1 if i -th ticket was for that movie
 - ▶ *Example*: Three movies M_1, M_2, M_3

$$\begin{array}{rcccccl} M_1 : & 0 & 0 & 0 & 1 & \dots \\ M_2 : & 1 & 0 & 0 & 0 & \dots \\ M_3 : & 0 & 1 & 1 & 0 & \dots \end{array} \quad (12)$$

First ticket to M_2 , second and third ticket to M_3 , fourth ticket to M_1

- ▶ Pick window of size N , where N is to reflect tickets to be recent

DECAYING WINDOWS: MOTIVATION

► *Possible solution:*

- *Example:* Three movies M_1, M_2, M_3

$$\begin{array}{rcccc} M_1 : & 0 & 0 & 0 & 1 & \dots \\ M_2 : & 1 & 0 & 0 & 0 & \dots \\ M_3 : & 0 & 1 & 1 & 0 & \dots \end{array} \quad (13)$$

First ticket to M_2 , second and third ticket to M_3 , fourth ticket to M_1

- Pick window of size N , where N is to reflect tickets to be recent
- Estimate number of ones in each stream
 - E.g. use Datar-Gionis-Indyk-Motwani (DGIM) algorithm
 - Estimates number of tickets sold for each movie
- Rank movies by the estimated counts

DECAYING WINDOWS: MOTIVATION

- ▶ *Possible solution, summary:*
 - ▶ One bit stream for each movie
 - ▶ i -th bit in a movie stream is 1 iff i -th ticket was for that movie
 - ▶ Count number of ones in each stream...
 - ▶ ... counts tickets for each movie
 - ▶ Rank movies by ticket counts
- ▶ Works for movies, because there only thousands of movies
- ▶ *Drawback:*
 - ▶ Does not work for items at Amazon or tweets per Twitter-user
 - ▶ 🗨️ too many items or users

DECAYING WINDOWS: MOTIVATION

- ▶ *Stream*: Movie tickets purchased all over the world
- ▶ *Goal*: Listing currently most “popular” movies
- ▶ *Alternative approach*:
 - ▶ Do not count ones in fixed-size window
 - ▶ Rather, compute “smooth aggregation” of *all* ones in stream
 - ▶ Smooth: use weights to rate stream elements in terms of recentness
 - ▶ The further back in the stream, the less weight given
 - ▶ *Example*: a_t most recent stream element

$$\begin{array}{rcccccc} \text{Stream :} & a_1 & a_2 & \cdots & a_{t-1} & a_t \\ \text{Weights :} & w_1 & w_2 & \cdots & w_{t-1} & w_t \end{array} \quad (14)$$

where

$$w_1 \leq w_2 \leq \dots \leq w_{t-1} \leq w_t$$

EXPONENTIALLY DECAYING WINDOW: DEFINITION

DEFINITION [EXPONENTIALLY DECAYING WINDOW]:

- ▶ Let a_1, a_2, \dots, a_t be a stream, with a_t most recent element
- ▶ Let c be small constant, e.g. $c \in [10^{-9}, 10^{-6}]$

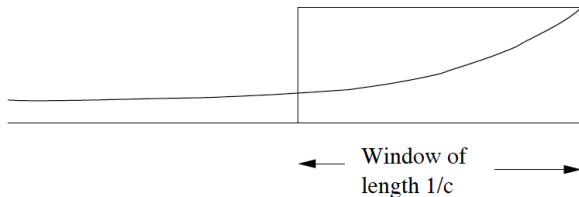
The *exponentially decaying window* for the stream is defined to be

$$\sum_{i=0}^{t-1} a_{t-i}(1-c)^i \quad (15)$$

Weight is $(1-c)^i$, it holds

$$(1-c)^0 \geq (1-c) \geq (1-c)^2 \geq (1-c)^3 \geq \dots \geq (1-c)^{t-1}$$

EXPONENTIALLY DECAYING WINDOW: DEFINITION



Decaying window and fixed-length window of equal weight

From mmds.org

- ▶ Decaying window puts weight $(1 - c)^i$ on $(t - i)$ -th element
- ▶ A window of length $1/c$ puts equal weight 1 on the first $1/c$ elements
- ▶ Both principles distribute the same weight to stream elements overall

UPDATING EXPONENTIALLY DECAYING WINDOWS

Upon arrival of a new element a_{t+1} , one updates the exponentially decaying window $\sum_{i=0}^{t-1} a_{t-i}(1-c)^i$ by

1. multiplying the current window by $(1-c)$, yielding

$$\sum_{i=0}^{t-1} a_{t-i}(1-c)^{i+1}$$

2. adding a_{t+1} , yielding

$$\sum_{i=0}^{t-1} a_{t-i}(1-c)^{i+1} + a_{t+1} = \sum_{i=0}^{(t+1)-1} a_{(t+1)-i}(1-c)^i$$

EXPONENTIALLY DECAYING WINDOWS: FINDING MOST POPULAR MOVIES

- ▶ *Most Popular Movies: Idea*
 - ▶ Have a bit stream for each movie, as before
 - ▶ Use e.g. $c = 10^{-9}$ (\approx sliding window of size $1/c = 10^9$)
 - ▶ On incoming movie ticket sale, update all decaying windows, as described above
 - ▶ First, multiply all decaying windows by $1 - c$
 - ▶ Add 1 for stream of the movie of the ticket; if there is no stream for that movie, create one
 - ▶ Do nothing (add 0) for all other streams
 - ▶ If any decaying window drops below threshold of $1/2$, drop window
 - ▶ Because the sum of all scores is $1/c$, there cannot be more than $2/c$ movies with score of $1/2$ or more
 - ▶ So, $2/c$ is limit on number of movies being tracked at any time
 - ▶ In practice, there should be much less movies counted
- ▶ *Therefore, one can apply the technique also for Amazon items and Twitter-users*

MATERIALS / OUTLOOK

- ▶ See *Mining of Massive Datasets*, chapter 4.6, 4.7
- ▶ As usual, see <http://www.mmids.org/> in general for further resources
- ▶ Next lecture: “Mining Data Streams III / Social Networks I”
 - ▶ See *Mining of Massive Datasets* 4.5; 10.1, 10.2, 10.3