

Experiment 7 : Implementing the DGIM Algorithm for Approximate Counting of Ones in a Data Stream

Date of Performance :

Date of Submission :

1. Statement of Problem

1.1 Aim:

To implement the DGIM (Datar-Gionis-Indyk-Motwani) algorithm in a programming language to efficiently estimate the number of 1's in the last k bits of a binary data stream using sublinear memory.

1.2 Objectives:

1. To understand the concept of streaming algorithms and their space-efficiency requirements.
2. To implement the DGIM algorithm for approximate counting of 1's in a binary stream using timestamps and bucket-based representation.
3. To simulate a real-time binary data stream and apply the DGIM algorithm for dynamic estimation.
4. To analyze the space complexity and error bounds of the DGIM algorithm compared to exact counting methods.
5. To evaluate the performance and accuracy of the DGIM implementation through experiments.

2. Theory

2.1 Data Stream Mining and Approximate Query Processing

(A) Data stream mining deals with extracting patterns, statistics, and trends from continuous, unbounded streams of data. These streams arrive in real-time and at high velocity, often making it infeasible to store and process the entire dataset. Traditional algorithms that assume random access to data or multiple passes over stored datasets fail in such settings. The DGIM algorithm is a streaming algorithm that works efficiently under the constraint of limited memory and single-pass access

(B) Approximate query processing (AQP) is essential when full precision is not mandatory but fast, scalable responses are. In the case of counting the number of 1's in a binary stream within a sliding window of the last k bits, exact counting would require storing the entire window. The DGIM algorithm addresses this challenge by allowing approximate counting with a guarantee on

the maximum error while using logarithmic space. It enables real-time estimation in streaming scenarios such as network monitoring, financial transactions, and sensor data collection.

2.2 Space-Efficient Counting Using DGIM Algorithm

(A) The DGIM algorithm uses a compact data structure based on buckets to summarize the positions of 1's in the stream. Each bucket stores two values: a timestamp or index of the most recent 1 it represents and the size of the bucket, which is always a power of two. Buckets of the same size are merged when more than two exist, ensuring minimal memory usage while preserving enough information for approximate counting.

(B) The key property of the DGIM algorithm is that it maintains at most two buckets of any given size. This strategy ensures that the total number of buckets remains small, and as a result, the space complexity is bounded by $O(\log^2 n)$, where n is the length of the stream. This makes the DGIM algorithm suitable for high-speed data streams where storing all bits is not possible. The merging process is done incrementally as new 1's are encountered, keeping the algorithm efficient and scalable

2.3 Real-Time Binary Stream Simulation

(A) Simulating a binary stream for DGIM involves generating a sequence of bits (0s and 1s) in real-time, mimicking the behavior of a live data source. Each bit is associated with a position index or a timestamp that helps identify whether it falls within the current sliding window. The simulation may use a random generator or predefined patterns to test the algorithm's behavior under various distributions of 1's and 0's.

(B) As the stream progresses, the DGIM algorithm responds dynamically to each incoming bit. When a 1 appears, a new bucket of size 1 is created, and older buckets are reviewed for merging based on the algorithm's rules. Buckets that fall outside the sliding window are discarded using their timestamps. This approach ensures that only relevant buckets are kept in memory and that the summary remains up to date with minimal overhead.

2.4 Query Processing and Error Bound Analysis

(A) When querying the number of 1's in the last k bits, the DGIM algorithm processes the maintained buckets in reverse order—from most recent to oldest. It adds the sizes of buckets whose timestamps fall entirely within the window. For the oldest bucket that partially overlaps the boundary of the window, only half of its count is considered in the estimate. This method allows efficient computation of the approximate count with only a few operations.

(B) The error introduced by the approximation is strictly bounded. The maximum error is at most the size of the partially included bucket, which can be no more than 50 percent of the actual value. Thus, the DGIM algorithm provides a guarantee that the answer will not differ from the true count by more than a known constant factor. This balance between accuracy and memory efficiency is what makes DGIM practical for real-time systems where quick, approximate answers are sufficient.

2.5 Evaluation of Accuracy and Performance

(A) To assess the effectiveness of the DGIM algorithm, its outputs are compared with those obtained from an exact counting method that stores all k bits and counts the number of 1's directly. The evaluation covers various parameters such as stream length, window size, and the density of 1's in the stream. Performance is measured in terms of memory usage, execution time, and deviation from the actual count.

(B) Experimental results are often visualized using graphs and tables to illustrate the trade-offs. For example, a plot of error percentage versus memory usage can demonstrate how increasing the window size affects accuracy. Tools like Matplotlib and Pandas can be used to perform the analysis. The results typically show that DGIM maintains acceptable accuracy while using significantly less memory, confirming its utility in high-volume, resource-constrained environments.

3 Implementation

3.1 Environment Setup

Algorithm: DGIM Algorithm for Approximate Counting of 1's in the Last k Bits of a Binary Data Stream

Algorithm Name: DGIM_Approximate_Count

Input Variables:

- `stream`: A binary data stream consisting of 0's and 1's.
- `k`: Window size (integer), i.e., the number of most recent bits to consider for the query.
- `t`: The current timestamp or index (automatically increments with the stream).

Output Variables:

- `approx_count`: Estimated number of 1's in the last k bits.

Variable Declarations:

- `Bucket`: A tuple (timestamp, size), where:
 - `timestamp`: integer representing the most recent occurrence of 1 in the bucket.
 - `size`: integer representing the number of 1's the bucket holds (always a power of 2).
- `Buckets`: A list of `Bucket`, representing all active buckets.
- `current_time`: integer, tracks the index of the latest bit in the stream.

Steps of the Algorithm:

Step 1: Initialization

1. Initialize an empty list `Buckets`.
2. Set `current_time = 0`.

Step 2: Stream Processing (for each bit)

1. For each bit in the `stream`:
 - Increment `current_time` by 1.
 - If the bit is 0:
 - Continue to the next bit.
 - If the bit is 1:
 - Create a new bucket: `(current_time, 1)`.
 - Insert this bucket at the beginning of the `Buckets` list.
 - Call `merge_buckets()` function to maintain at most 2 buckets per size.

Step 3: Merging Buckets (Function: `merge_buckets`)

1. Traverse the `Buckets` list from newest to oldest.
2. For each group of 3 consecutive buckets with the same size:
 - Merge the two oldest:
 - New size = sum of their sizes ($2 \times s$).
 - New timestamp = timestamp of the most recent among the two.
 - Replace the two oldest with the merged bucket.
3. Repeat this until no size has more than two buckets.

Step 4: Query Approximate Count (Function: `estimate_ones_in_last_k(k)`)

1. Set `threshold_time = current_time - k`.
2. Initialize `approx_count = 0`.
3. Traverse `Buckets` from newest to oldest:
 - If bucket's timestamp > `threshold_time`:
 - Add `bucket.size` to `approx_count`.
 - If bucket's timestamp <= `threshold_time`:
 - Add `bucket.size / 2` to `approx_count` (partial bucket considered).
 - Break the loop.
4. Return `approx_count`.

3.2 Explanation of Statement of Problem

In many real-world applications such as network traffic monitoring, clickstream analysis, and financial data processing, data arrives continuously in the form of streams. Storing the entire stream is impractical due to its unbounded size, and performing exact computations on such large streams is inefficient.

A common task is to determine the number of 1s (bits set to one) in the most recent N elements of a binary data stream (sliding window model). A naive approach would require storing all N bits and performing exact counting, which demands $O(N)$ space and time, making it infeasible for high-speed or infinite streams.

The DGIM (Datar–Gionis–Indyk–Motwani) Algorithm provides an efficient way to approximately count the number of 1s in the most recent N bits of a data stream. It reduces memory usage from $O(N)$ to $O(\log^2 N)$ while still guaranteeing a bounded error in the result.

Therefore, the problem is to implement the DGIM algorithm for approximate counting of 1s in a binary data stream, ensuring correctness, efficiency, and error bounds, while avoiding the need to store the entire stream.

3.3 Execute

Program:

```
class DGIMCounter:
    def __init__(self):
        self.bucket_list = [] # Stores (timestamp, size)
        self.time = 0

    def _compress_buckets(self):
        idx = 0
        while idx < len(self.bucket_list) - 2:
            if (self.bucket_list[idx][1] == self.bucket_list[idx+1][1] ==
                self.bucket_list[idx+2][1]):
                merged_time = self.bucket_list[idx+1][0]
                merged_size = self.bucket_list[idx+1][1] * 2
                del self.bucket_list[idx+1]
                del self.bucket_list[idx+1]
                self.bucket_list.insert(idx+1, (merged_time, merged_size))
```

```

        else:
            idx += 1

def add_bit(self, bit):
    self.time += 1
    if bit == 1:
        self.bucket_list.insert(0, (self.time, 1))
        self._compress_buckets()

def estimate_count(self, k):
    boundary = self.time - k
    total = 0
    for ts, size in self.bucket_list:
        if ts > boundary:
            total += size
        else:
            total += size // 2
            break
    return total

```

3.4 Output

- **Case1: Exact Match (Absolute Error = 0)**

Input:

```

Binary Stream: [1, 0, 1, 1, 0, 0,
1, 0, 1, 1, 1, 1, 1, 1, 1 0]
Window Size (k): 6

```

Output:

```

--- DGIM Approximation ---
Estimated 1's in last 6 bits: 3
Actual 1's in last 6 bits: 3
Absolute Error: 0

```

In this case, all buckets representing 1's are fully contained within the last k bits. Since no bucket is partially outside the window, the DGIM estimate exactly matches the actual count, resulting in zero error

- **Case2: Approximate Result (Absolute Error > 0)**

Input:

Binary Stream: [1, 0, 1, 1, 0, 0, 1, 1, 0, 1, 1, 1, 1, 0, 1, 1, 0, 1, 1, 0]

Output:

```
-- DGIM Approximation --
Estimated 1's in last 6 bits: 4
Actual 1's in last 6 bits:    3
Absolute Error: 1
```

In this case, the oldest bucket within the last k bits is **partially inside the window**. DGIM cannot know the exact number of 1's in that portion, so it assumes **half of the bucket's size** falls inside. This leads to a small overestimation (4 instead of the actual 3), which is expected behavior and is within the algorithm's guaranteed error bound of 50%.

Key Notes from Theory to Code

Concept	Implementation
Bit stream arrival	<code>process_bit()</code> method simulates live arrival
Buckets	Stored as <code>(timestamp, size)</code>
Bucket merging	<code>merge_buckets()</code> ensures 2-bucket rule
Error bound	Worst-case overcount $\leq \frac{1}{2}$ bucket size
Space complexity	$O(\log_2 n)$ active buckets
Query time	Single traversal of <code>buckets</code> list

4 Conclusion:

The DGIM algorithm efficiently estimates the number of 1's in the last k bits of a binary data stream using limited memory. It summarizes the stream into logarithmically sized buckets, allowing approximate counting without storing the entire stream. In Case 1, the algorithm gives an exact count, showing it can be perfectly accurate under certain conditions. In Case 2, the estimate differs by 1 due to partial bucket approximation, illustrating the trade-off between accuracy and memory usage. The error is always small and predictable, making the algorithm reliable for streaming data. DGIM's memory requirements grow logarithmically with the stream size, which is far less than storing all bits. The algorithm updates continuously as new bits arrive, supporting real-time monitoring. Its bucket-based design is simple to implement in software and hardware. By balancing minor errors with huge memory savings, DGIM is highly efficient. Overall, this experiment demonstrates that DGIM is both practical and effective for real-world data stream applications.

5. References:

- Datar, M., Gionis, A., Indyk, P., & Motwani, R. (2002). *Maintaining stream statistics over sliding windows*. SIAM Journal on Computing, 31(6), 1794–1813.
- Leskovec, J., Rajaraman, A., & Ullman, J. D. (2014). *Mining of Massive Datasets*. Cambridge University Press.
- Maheshwari, A. (2019). *Stream statistics over sliding windows*. Carleton University.
- Raj, V. (2020). *DGIM Algorithm Implementation in C*. GitHub Repository.

Sign and Remark:

R1 (4 Marks)	R2 (4 Marks)	R3 (4 Marks)	R4 (3 Marks)	Total Marks (15 Marks)	Signature