

SCE17-0037:
**Experimental Analysis of the
Application of the gSketch
Partitioning Method onto the
gMatrix Graph-Stream Sketch**



Eric Leonardo Lim
U1420158D

Supervisor: Ast/P Arijit Khan

School of Computer Science and Engineering
Nanyang Technological University

Submitted in partial fulfilment of the requirements for the degree of
Bachelor of Engineering (Computer Science)

AY2017/2018

Acknowledgements

I would like to express my gratitude to everyone who gave me the possibility to complete this report. A special thanks to my final year project supervisor, Asst/Prof Arijit Khan, who has provided me with relevant guidance on the project.

Abstract

This report presents a final year project that is about an experimental analysis of applying the gSketch partitioning method onto the gMatrix graph-stream sketch. The report first introduces how the gSketch partitioning method can be applied onto the gMatrix sketch and proposes optimizations for the method, and then analyzes how the gSketch partitioning method changes how gMatrix answers various query types, such as edge frequency, heavy-hitter edges, and node aggregate-frequency queries, and how the performance and probabilistic accuracy guarantees change, and after that, shows experimental results with metrics that each evaluates differently how partitioning affects gMatrix's accuracy for answering the different query types on up to three different graph-stream datasets. Finally, the report concludes that the gSketch partitioning method successfully improves the accuracy of gMatrix in query types such as edge frequency estimation and source-node aggregate-frequency estimation, although fails to bring the same improvements onto the destination-node aggregate-frequency estimation and heavy-hitter edge queries.

Table of contents

List of figures	ix
1 Introduction	1
1.1 Background	1
1.2 Objectives	2
2 Design and Analysis	3
2.1 Preliminaries	3
2.1.1 Count-Min Sketch	3
2.1.2 gMatrix Sketch	3
2.1.3 gSketch Partitioning Method	4
2.2 gSketch Partitioning Method for gMatrix	5
2.2.1 Edge Frequency Query Estimation	7
2.2.2 Heavy-hitter Edge Query Estimation	9
2.2.3 Node Aggregate-Frequency Query Estimation	11
3 Experiments	15
3.1 Implementation	15
3.2 Datasets	16
3.3 Metrics	16
3.3.1 Edge Frequency & Node Aggregate-Frequency Estimation	16
3.3.2 Heavy Hitter Edge Estimation	17
3.4 Results	18
3.4.1 Edge Frequency Estimation	18
3.4.2 Heavy Hitter Edge Estimation	28
3.4.3 Node Aggregate-Frequency Estimation	29
4 Conclusions	33
4.1 Conclusion	33

Table of contents

4.2 Possible Improvements	33
References	35

List of figures

3.1	Observed error on Friendster dataset considered over 1 million random edges	18
3.2	Observed error on Friendster dataset considered over top-500 highest frequency edges	19
3.3	Average relative error on Friendster dataset considered over 1 million random edges	19
3.4	Average relative error on Friendster dataset considered over top-500 highest frequency edges	20
3.5	Effective queries percentage on Friendster dataset considered over 1 million random edges	20
3.6	Effective queries percentage on Friendster dataset considered over top-500 highest frequency edges	21
3.7	Observed error of Twitter dataset considered over 1 million random edges	21
3.8	Observed error of Twitter dataset considered over top-500 highest frequency edges	22
3.9	Average relative error on Twitter dataset considered over 1 million random edges	22
3.10	Average relative error on Twitter dataset considered over top-500 highest frequency edges	23
3.11	Effective queries percentage on Twitter dataset considered over 1 million random edges	23
3.12	Effective queries percentage on Twitter dataset considered over top-500 highest frequency edges	24
3.13	Observed error of IP-Trace dataset considered over 1 million random edges	24
3.14	Observed error of IP-Trace dataset considered over top-500 highest frequency edges	25

List of figures

3.15	Average relative error on IP-Trace dataset considered over 1 million random edges	25
3.16	Average relative error on IP-Trace dataset considered over top-500 highest frequency edges	26
3.17	Effective queries percentage on IP-Trace dataset considered over 1 million random edges	26
3.18	Effective queries percentage on IP-Trace dataset considered over top-500 highest frequency edges	27
3.19	False positive rate of heavy hitter edge estimation on IP-Trace dataset with respect to varying frequency threshold percentages of the total stream frequency	28
3.20	False positive rate of heavy hitter edge estimation on Friendster dataset with respect to varying frequency threshold percentages of the total stream frequency	28
3.21	Observed error of source-node aggregate-frequency queries on top-500 node aggregate-frequencies of the Friendster dataset	29
3.22	Observed error of destination-node aggregate-frequency queries on top-500 node aggregate-frequencies of the Friendster dataset	30
3.23	Observed error of source-node aggregate-frequency queries on top-500 node aggregate-frequencies of the Twitter dataset	30
3.24	Observed error of destination-node aggregate-frequency queries on top-500 node aggregate-frequencies of the Twitter dataset	30
3.25	Observed error of source-node aggregate-frequency queries on top-500 node aggregate-frequencies of the IP-Trace dataset	31
3.26	Observed error of destination-node aggregate-frequency queries on top-500 node aggregate-frequencies of the IP-Trace dataset	31

Chapter 1

Introduction

1.1 Background

Graphs are naturally used to model data that consists of entities and the relationships between them, with entities represented as nodes and the relationships between a pair of entities represented as edges. Common applications include modelling communication networks, social networks, and the Web. In such applications, graph data is often available as a massive rapid input stream of edges, so computing exact properties of the graph is often not computationally feasible. Thus, sketches, which summarize general data streams, are often used to approximate queries with a particular guaranteed accuracy.

Previous work [1] has introduced a way to utilize the Count-min sketch [2], a general data-stream sketch for approximating frequency counts, to summarize graph streams. Furthermore, it [1] has introduced the gSketch, which is a partitioning method that exploits typical local structural properties of real world graph streams by building multiple localized Count-min sketches [2] for answering graph stream edge-frequency estimation. It has been shown that the partitioning method is able to improve overall accuracy of the Count-min sketch when tested on several large graph datasets.

A more recent work [3] has introduced the gMatrix, which is a sketch that generalizes the Count-min sketch for graph streams and specializes only on graph streams. The gMatrix has been shown to be able to handle various queries involving the structural properties of the underlying graph, which is an advantage over the Count-min sketch [1].

1.2 Objectives

The gSketch partitioning method improves the Count-Min sketch significantly. An idea is to implement the method into the gMatrix [3] to improve its accuracy. However, there are no previous works that have shown if it really is the case. In addition, since the main feature of the gMatrix sketch is its support for various different query types, we want to know whether partitioning can bring improvements onto query types other than edge frequency queries.

As there are no previous works that have attempted to analyze the effects of the application of the gSketch partitioning method onto the gMatrix, the project is conducted with the objective of analyzing how the gSketch[1] partitioning method affects the gMatrix[3] with the goal of improving its accuracy on answering various query types.

Chapter 2

Design and Analysis

2.1 Preliminaries

2.1.1 Count-Min Sketch

The Count-Min [2] sketch is a data-stream sketch that is used to approximate frequency counts of items in a data stream.

It consists of a 2D array. In a $h \times w$ sized Count-Min sketch, it consists of w many hash functions that each maps onto $[0, h - 1]$.

To store the frequency of an item x , it first maps x onto w different values $g_k(x) \in [0, h - 1]$ using each hash function $g_k, k \in \{1..w\}$. Then, the entry at the $g_k(x)^{th}$ row, k^{th} column is incremented by the frequency of x .

Let the entry of the sketch at the i^{th} row, j^{th} column be denoted by $f(i, j)$. Querying the frequency estimate of an item x is accomplished by finding $\min\{f(g_k(x), k) | k \in \{1..w\}\}$.

2.1.2 gMatrix Sketch

The gMatrix [3] is a sketch that is similar to the Count-Min sketch, except that it consists of a 3D array instead of a 2D one and specializes on graph streams.

In a $h \times h \times w$ sized gMatrix, to store the frequency of an edge (i, j) , it first maps i and j onto w different values $g_k(i)$ and $g_k(j)$ using each hash function $g_k, k \in \{1..w\}$. Then, the entry at the $g_k(i)^{th}$ row, $g_k(j)^{th}$ column, and k^{th} layer is incremented by the frequency of the edge (i, j) .

Let the entry of the sketch at the i^{th} row, j^{th} column, and k^{th} layer be denoted by $f(i, j, k)$. Querying the frequency estimate of an edge (i, j) is accomplished by finding $\min\{f(g_k(i), g_k(j), k) | k \in \{1..w\}\}$.

2.1.3 gSketch Partitioning Method

Algorithm 1 Sketch-Partitioning-On-CountMin [1] (Data Sample: D)

```

1: Create a root node  $S$  of the partitioning tree as an active node;
2:  $S.width \leftarrow h$ ;
3:  $S.depth \leftarrow w$ ;
4: Create an empty list  $L$  containing  $S$  only;
5: while  $L \neq \emptyset$  do
6:   Partition active node  $S \in L$  based on  $D$  into  $S_1, S_2$  by minimizing  $E$  in equation
   2.1;
7:    $S_1.width = S_2.width = \frac{S.width}{2}$ ;
8:    $L \leftarrow L \setminus \{S\}$ ;
9:   if  $(S_1.width \geq w_0)$  and  $(\sum_{m \in S_1} \tilde{d}(m) > C \cdot S_1.width)$  then
10:     $L \leftarrow L \cup S_1$ ;
11:   else
12:    Construct the localized sketch  $S_1$ ;
13:   if  $(S_2.width \geq w_0)$  and  $(\sum_{m \in S_2} \tilde{d}(m) > C \cdot S_2.width)$  then
14:     $L \leftarrow L \cup S_2$ ;
15:   else
16:    Construct the localized sketch  $S_2$ ;

```

The gSketch [1] is a partitioning method performed onto a graph-stream sketch prior to being populated with the stream. The partitioning method is based on the source vertex set of a sample stream. The method, as shown in Algorithm 1 [1], is recursive and on each step of the recursion, it first sorts the sample in order of non-decreasing average edge frequency, which is estimated by the ratio of vertex aggregate-frequency to its degree $\frac{f_v(m)}{d(m)}$, and then selects a pivot on which the sorted list is to be partitioned into two by minimizing an objective function,

$$E = \sum_{m \in S_1} \frac{\tilde{d}(m) \cdot \tilde{F}(S_1)}{\tilde{f}_v(m)/\tilde{d}(m)} + \sum_{m \in S_2} \frac{\tilde{d}(m) \cdot \tilde{F}(S_2)}{\tilde{f}_v(m)/\tilde{d}(m)} \quad (2.1)$$

, where S_1 and S_2 are the two produced partitions based on the selected pivot on the step of the recursion and $\tilde{F}(S_i)$ is its total frequency, and that when minimized, the overall relative error of the whole sketch is minimized as well.

2.2 gSketch Partitioning Method for gMatrix

As shown in Algorithm 1, there are two conditions in which if any one of the two is met, the recursion is immediately terminated and starts to build the local partition S :

1. $S.width < w_0$ as the sketch is small enough.
2. $\sum_{m \in S} \tilde{d}(m) \leq C \cdot S.width$, where C is a constant, as it implies that the probability for any collisions to occur in any cell of S can be bounded above by C .

2.2 gSketch Partitioning Method for gMatrix

The gSketch partitioning[1] is shown to improve Count-Min's[2] accuracy. The partitioning works by grouping edges with similar frequencies in the same partition to improve sketch accuracy and requires an available data sample.

The same partitioning method can also be applied onto the gMatrix, as shown in Algorithm 2, by changing the width $S.width$ of the Count-Min sketch with the number of rows of the gMatrix $S.rows$.

Algorithm 2 Sketch-Partitioning-On-gMatrix (Data Sample: D)

```

1: Create a root node  $S$  of the partitioning tree as an active node;
2:  $S.rows \leftarrow h$ ;
3:  $S.cols \leftarrow h$ ;
4:  $S.depth \leftarrow w$ ;
5: Create an empty list  $L$  containing  $S$  only;
6: while  $L \neq \emptyset$  do
7:   Partition active node  $S \in L$  based on  $D$  into  $S_1, S_2$  by minimizing  $E$  in equation
   2.1;
8:    $S_1.rows = S_2.rows = \frac{S.rows}{2}$ ;
9:    $L \leftarrow L \setminus \{S\}$ ;
10:  if  $(S_1.rows \geq w_0)$  and  $(\sum_{m \in S_1} \tilde{d}(m) > C \cdot S_1.rows)$  then
11:     $L \leftarrow L \cup S_1$ ;
12:  else
13:    Construct the localized sketch  $S_1$ ;
14:  if  $(S_2.rows \geq w_0)$  and  $(\sum_{m \in S_2} \tilde{d}(m) > C \cdot S_2.rows)$  then
15:     $L \leftarrow L \cup S_2$ ;
16:  else
17:    Construct the localized sketch  $S_2$ ;

```

We provide a proof that the terminating condition $\sum_{m \in S} \tilde{d}(m) \leq C \cdot S.rows$ still guarantees that the probability for any collisions to occur in the sketch partition is bounded above by C .

Design and Analysis

Theorem 1. *Let S be a gMatrix partition and $0 < C < 1$ be a constant. Then, the condition $\sum_{m \in S} \tilde{d}(m) \leq C \cdot S.rows$ implies that the probability for any collisions to occur in S is bounded above by C .*

Proof. There are two cases of spurious edges that may collide with an edge (i, j) in S . The first case of such edges are those for which both end-points are neither i nor j , and the probability of collision is $\frac{1}{S.rows \cdot S.cols}$. The second case of such edges are those for which either the source node is i too or the destination node is j too, and the probability of collision is $\frac{1}{S.cols}$ and $\frac{1}{S.rows}$ consecutively. Since the partitioning divides the number of rows into half on each step of the recursion, on an $h \times h \times w$ gMatrix, it always holds that $S.rows \leq S.cols$, and so the probability of collision is at most $\frac{1}{S.rows}$. Since there are $\sum_{m \in S} \tilde{d}(m)$ distinct edges in S , the probability for any collisions to occur is at most $\frac{\sum_{m \in S} \tilde{d}(m)}{S.rows}$. Therefore, $\sum_{m \in S} \tilde{d}(m) \leq C \cdot S.rows$ guarantees that the probability for any collisions to occur in the sketch is bounded above by C . \square

After partitioning, each of the resulting partitions is associated with a set of source vertices. Each source vertex is then only associated with exactly one partition. Edge (u, v) is stored in partition p if and only if u is in the set of source vertices associated with the partition p . An outlier partition is then needed to be reserved to store frequencies of edges in the dataset whose source node is not associated with any of the resulting partitions of the sketch partitioning algorithm.

We propose two new optimizations to the partitioning method. The first is that since the partitioning works by grouping edges with similar average frequencies together estimated by the ratio of source-vertex aggregate-frequency to its degree $\frac{\tilde{f}_v(m)}{d(m)}$ and that edges with a common source-node m are assumed to have similar frequencies, we can filter ones for which the assumption does not hold. In particular, we can compute the statistical variances of the frequencies of edges from the data sample with common source nodes and filter edges with source nodes for which the variance of the frequencies is exceeding a certain threshold.

Another observation is that the outlier partition size to be used for the sketch can be estimated by computing the outlier ratio from the data sample. We do so by first splitting the data sample into two, and then after selecting the suitable source nodes from the first split, the ratio of outliers in the second split is computed, which is essentially the ratio of the number of source nodes in the split not found among the previously-selected source nodes to the size of the split.

Next, we observe how gMatrix answers various query types after being partitioned and how partitioning affects gMatrix's accuracy.

2.2.1 Edge Frequency Query Estimation

Since there is only one partition associated with each source vertex, to retrieve the estimated frequency for an arbitrary edge (i, j) :

1. Find partition p that is responsible for i
2. Find $\min\{V_P(g_k(i), g_k(j), k) | \forall k \in \{1..w\}\}$, where V_P is the value of the cell in the partition p .

The number of operations taken by the first step depends on implementation. If a hash-table that maps each source-vertex onto its allocated partition is to be built during the sketch partitioning algorithm, the first step takes $\mathcal{O}(1)$ operations in average. For worst-case time complexity analysis, if a balanced binary search tree is to be used instead of the hash-table, the first step takes $\mathcal{O}(\lg(|V^+|))$ operations, where V^+ is the set of source vertices/nodes from the data sample.

The second step takes $\mathcal{O}(w)$ operations.

Thus, the worst-case time complexity is $\mathcal{O}(\lg(|V^+|) + w)$.

Since obtaining the frequency estimate of an arbitrary edge only depends on the partition containing the edge, we observe how each partition answers the query.

Let S be a gMatrix of size $h \times h \times w$. Suppose the partitioning step has been performed already. Note that any of the produced partitions will have number of rows equal to $\frac{h}{2^d}$ where d is the depth of the recursion in which the partition is built. Let p be an arbitrary partition. The size of p is $\frac{h}{2^d} \times h \times w$.

Theorem 2. *Let p be a partition built at depth d of the sketch partitioning algorithm. Let (i, j) be an edge in p , $Q(i, j)$ be its actual frequency, and $\overline{Q(i, j)}$ be its estimated frequency according to partition p . Let L_p be the total frequency of edges received so far in the arbitrary partition, i.e. total frequency of edges (u, v) such that u is associated with p . Let $A_p(j)$ be the sum of the frequencies of edges (u, v) on p such that $v = j$. Let $B_p(i)$ be the sum of the frequencies of edges (u, v) on p such that $u = i$. Let $\epsilon \in (0, 1)$ be a very small fraction. Then,*

$$P(Q(i, j) \leq \overline{Q(i, j)}) \leq Q(i, j) + L_p \cdot \epsilon + A_p(j) \cdot h \cdot \epsilon + B_p(i) \cdot h \cdot \epsilon \geq 1 - \left(\frac{2^{d+1} + 1}{h^2 \cdot \epsilon}\right)^w$$

Proof. Note that $\overline{Q(i, j)}$ is essentially $Q(i, j)$ added with the frequencies of spurious edges mapped into the same cell $(g_k(i), g_k(j), k)$ in the partition, so $P(Q(i, j) \leq \overline{Q(i, j)}) = 1$. We remain to show that

$$P(\overline{Q(i, j)}) \leq Q(i, j) + L_p \cdot \epsilon + A_p(j) \cdot h \cdot \epsilon + B_p(i) \cdot h \cdot \epsilon \geq 1 - \left(\frac{2^{d+1} + 1}{h^2 \cdot \epsilon}\right)^w$$

There are three possible cases for the spurious edges.

The first possible case of spurious edges (u, v) are those for which $u \neq i$ and $v \neq j$. Any of such edges are equally likely to be mapped into any cell in p , and the probability to be mapped into any particular cell is $\frac{1}{h \cdot \frac{h}{2^d}} = \frac{2^d}{h^2}$. Let X_k be a random variable that represents the number of such spurious edges that are mapped into $(g_k(i), g_k(j), k)$. Therefore, $E[X_k] = \frac{2^d L_p}{h^2}$. Then, by Markov's inequality,

$$P(X_k \geq L_p \cdot \epsilon) \leq \frac{E[X_k]}{L_p \cdot \epsilon} = \frac{2^d}{h^2 \epsilon} \quad (2.2)$$

The second possible case of spurious edges (u, v) are those for which $u \neq i$ but $v = j$. The probability for any of such edges to be mapped into $(g_k(i), g_k(j), k)$ in p is $\frac{1}{h} = \frac{2^d}{h}$. Let Y_k be a random variable representing the number of such spurious edges. Therefore, $E[Y_k] = \frac{2^d A_p(j)}{h}$. By Markov's inequality,

$$P(Y_k \geq h A_p(j) \cdot \epsilon) \leq \frac{E[Y_k]}{h A_p(j) \cdot \epsilon} = \frac{2^d}{h^2 \epsilon} \quad (2.3)$$

The third possible case of spurious edges (u, v) are those for which $u = i$ but $v \neq j$. The probability for any of such edges to be mapped into $(g_k(i), g_k(j), k)$ in p is $\frac{1}{h}$. Let Z_k be a random variable representing the number of such spurious edges. Therefore, $E[Z_k] = \frac{B_p(i)}{h}$. By Markov's inequality,

$$P(Z_k \geq h B_p(i) \cdot \epsilon) \leq \frac{E[Z_k]}{h B_p(i) \cdot \epsilon} = \frac{1}{h^2 \epsilon} \quad (2.4)$$

Combining inequations 2.2, 2.3, 2.4,

$$P(X_k + Y_k + Z_k \geq L_p \cdot \epsilon + A_p(j) \cdot h \cdot \epsilon + B_p(i) \cdot h \cdot \epsilon) \leq \frac{2^{d+1} + 1}{h^2 \cdot \epsilon} \quad (2.5)$$

Therefore,

$$\begin{aligned}
P(\overline{Q(i, j)}) &\leq Q(i, j) + L_p \cdot \epsilon + A_p(i) \cdot h \cdot \epsilon + B_p(i) \cdot h \cdot \epsilon \\
&= 1 - \prod_{k=1}^w P(X_k + Y_k + Z_k \geq L_p \cdot \epsilon + A_p(j) \cdot h \cdot \epsilon + B_p(i) \cdot h \cdot \epsilon) \\
&\geq 1 - \left(\frac{2^{d+1} + 1}{h^2 \cdot \epsilon}\right)^w
\end{aligned} \tag{2.6}$$

□

Remarks. *The probability of the guarantee gets lower as the depth d of the recursion in which a partition is built gets deeper, but in ideal partitioning, the guarantee of the estimate of the partition itself gets better as L_p , $A_p(i)$, and $B_p(i)$ are reduced.*

2.2.2 Heavy-hitter Edge Query Estimation

The query asks for retrieving the set of edges in the graph stream with frequencies at least F . The query is accomplished by the following steps:

1. obtain the set of edges with frequencies at least F from each partition in the gMatrix, which can be accomplished as described in [3]
2. compute the union of each of the obtained sets.

Without partitioning, we only need to perform step 1 once, but with partitioning, we need to perform it n times, where n is the number of partitions, so it is already at least n times slower than without partitioning in the worst case. Thus, the first step takes $\mathcal{O}(nh^2w + 2n \lfloor P/h \rfloor \log P \prod_{k=1}^w c_k)$ operations [3], where c_k is the number of cells in the k^{th} layer whose entries are at least F , and P is a gMatrix parameter for generating hash functions and is a prime number larger than the maximum number of nodes in the graph stream.

For the second step, if the partitioning algorithm produces only one partition, there are no union operations performed. Otherwise, the second step takes $\Omega(|H|)$ operations, where H denotes the set of true heavy-hitter edges in the graph stream.

As in [3], the resulting set of heavy-hitter edges may consist of edges that are not actually heavy-hitter edges, but the set of all true heavy-hitter edges is always a subset, i.e. there are no missing true heavy-hitter edges in the resulting set.

We observe the probability that an edge (i, j) is correctly classified as a heavy-hitter edge in each partition.

Theorem 3. *Let p be a partition built at depth d of the sketch partitioning algorithm. Let (i, j) be an edge in p . Let $Q(i, j)$ be its actual frequency, and $\overline{Q(i, j)}$ be its estimated frequency according to partition p . Let L_p be the total frequency of edges received so far in the arbitrary partition, i.e. total frequency of edges (u, v) such that u is associated with p . Let $A_p(j)$ be the sum of the frequencies of edges (u, v) on p such that $v = j$. Let $B_p(i)$ be the sum of the frequencies of edges (u, v) on p such that $u = i$. Let $\epsilon \in (0, 1)$ be a very small fraction. Then,*

$$P(Q(i, j) \geq F) \geq 1 - \left(\frac{3 \cdot 2^d (L_p + h \cdot A_p(j)) + 3 \cdot h \cdot B_p(i)}{h^2 \cdot (\overline{Q(i, j)} - F)} \right)^w$$

Proof. Note that if the number of spurious edges is at most $(\overline{Q(i, j)} - F)$, the true frequency $Q(i, j)$ is at least F , so by obtaining a lower bound for the probability of the former, we obtain a lower bound for the probability of the latter as well.

To obtain a lower bound for the probability that the number of spurious edges is at most $(\overline{Q(i, j)} - F)$, we consider all possible cases of spurious edges.

The first possible case of spurious edges (u, v) are those for which $u \neq i$ and $v \neq j$. Any of such edges are equally likely to be mapped into any cell in p , and the probability to be mapped into any particular cell is $\frac{1}{h \cdot 2^d} = \frac{2^d}{h^2}$. Let X_k be a random variable that represents the number of such spurious edges that are mapped into $(g_k(i), g_k(j), k)$. Therefore, $E[X_k] = \frac{2^d L_p}{h^2}$. Then, by Markov's inequality,

$$P(X_k \geq \frac{\overline{Q(i, j)} - F}{3}) \leq \frac{E[X_k]}{\frac{\overline{Q(i, j)} - F}{3}} = \frac{3 \cdot 2^d L_p}{h^2 (\overline{Q(i, j)} - F)} \quad (2.7)$$

The second possible case of spurious edges (u, v) are those for which $u \neq i$ but $v = j$. The probability for any of such edges to be mapped into $(g_k(i), g_k(j), k)$ in p is $\frac{1}{h} = \frac{2^d}{h}$. Let Y_k be a random variable representing the number of such spurious edges. Therefore, $E[Y_k] = \frac{2^d A_p(j)}{h}$. By Markov's inequality,

$$P(Y_k \geq \frac{\overline{Q(i, j)} - F}{3}) \leq \frac{E[Y_k]}{\frac{\overline{Q(i, j)} - F}{3}} = \frac{3 \cdot 2^d A_p(j)}{h (\overline{Q(i, j)} - F)} \quad (2.8)$$

The third possible case of spurious edges (u, v) are those for which $u = i$ but $v \neq j$. The probability for any of such edges to be mapped into $(g_k(i), g_k(j), k)$ in p is $\frac{1}{h}$. Let Z_k be a random variable representing the number of such spurious edges. Therefore, $E[Z_k] = \frac{B_p(i)}{h}$. By Markov's inequality,

$$P(Z_k \geq \frac{\overline{Q(i,j)} - F}{3}) \leq \frac{E[Z_k]}{\frac{\overline{Q(i,j)} - F}{3}} = \frac{3 \cdot B_p(i)}{h(\overline{Q(i,j)} - F)} \quad (2.9)$$

Combining inequations 2.7, 2.8, 2.9,

$$P(X_k + Y_k + Z_k \geq \overline{Q(i,j)} - F) \leq \frac{3 \cdot 2^d(L_p + h \cdot A_p(j)) + 3 \cdot h \cdot B_p(i)}{h^2 \cdot (\overline{Q(i,j)} - F)} \quad (2.10)$$

Therefore,

$$\begin{aligned} P(Q(i,j) \geq F) \\ &= 1 - \prod_{k=1}^w P(X_k + Y_k + Z_k \geq \overline{Q(i,j)} - F) \\ &\geq 1 - \left(\frac{3 \cdot 2^d(L_p + h \cdot A_p(j)) + 3 \cdot h \cdot B_p(i)}{h^2 \cdot (\overline{Q(i,j)} - F)} \right)^w \end{aligned} \quad (2.11)$$

□

2.2.3 Node Aggregate-Frequency Query Estimation

The queries ask for the sum of frequencies of all edges incoming to / outgoing from a node i . Due to partitioning based on source nodes, source-node aggregate-frequency queries and destination-node aggregate-frequency queries are answered differently.

Source-Node Aggregate-Frequency Query Estimation

The task of finding the aggregate frequency of a source node i is broken down to the following steps:

1. Finding the associated partition p of i .
2. Finding the aggregate frequency of the source node i at the partition p , which can be accomplished as described in [3].

The number of operations taken by the first step depends on implementation. If a hash-table that maps each source-vertex onto its allocated partition is to be built during the sketch partitioning algorithm, the first step takes $\mathcal{O}(1)$ operations in average. For worst-case time complexity analysis, if a balanced binary search tree is to be used instead of the hash-table, the first step takes $\mathcal{O}(\lg(|V^+|))$ operations, where V^+ is the set of source vertices/nodes from the data sample.

Design and Analysis

The second one takes $\mathcal{O}(hw)$ operations, so the query can be answered in $\mathcal{O}(hw)$ operations.

Thus, the overall time complexity is $\mathcal{O}(\lg(|V^+|) + hw)$.

Next, we observe how partitioning affects the probabilistic accuracy guarantee.

Theorem 4. *Let p be a partition built at depth d of the sketch partitioning algorithm. Let (i, j) be an edge in p . Let $Q_{agg}^+(i)$ be the true aggregate-frequency of source-node i , and $\overline{Q_{agg}^+(i)}$ be its estimated aggregate-frequency. Let L_p be the total frequency of edges received so far on p . Let $\epsilon \in (0, 1)$ be a very small fraction. Then,*

$$P(Q_{agg}^+(i) \leq \overline{Q_{agg}^+(i)} \leq Q_{agg}^+(i) + L_p \cdot \epsilon) \geq 1 - \left(\frac{2^d}{h \cdot \epsilon}\right)^w$$

Proof. We note that $P(Q_{agg}^+(i) \leq \overline{Q_{agg}^+(i)}) = 1$ since $\overline{Q_{agg}^+(i)}$ is always an overestimate. We remain to show that

$$P(\overline{Q_{agg}^+(i)} \leq Q_{agg}^+(i) + L_p \cdot \epsilon) \geq 1 - \left(\frac{2^d}{h \cdot \epsilon}\right)^w$$

In any arbitrary partition p , the probability for a spurious edge (u, v) , where $u \neq i$ and u is associated with p , to be mapped onto any of the columns at the $g_k(i)^{th}$ row, k^{th} layer of p is $\frac{1}{h} = \frac{2^d}{h}$, so the expected number of spurious edges is $\frac{2^d L_p}{h}$. Let R_k be the random variable that represents the number of such spurious edges for the k^{th} hash function. By Markov's inequality,

$$P(R_k \geq L_p \epsilon) \leq \frac{E[R_k]}{L_p \epsilon} = \frac{2^d}{h \epsilon} \quad (2.12)$$

Therefore,

$$\begin{aligned} P(\overline{Q_{agg}^+(i)} \leq Q_{agg}^+(i) + L_p \cdot \epsilon) \\ &= 1 - \prod_{k=1}^w P(R_k \geq L_p \cdot \epsilon) \\ &\geq 1 - \left(\frac{2^d}{h \cdot \epsilon}\right)^w \end{aligned} \quad (2.13)$$

□

Remarks. *As the depth d of the recursion in which the partition is built increases, the probability of the accuracy guarantee decreases, but in ideal partitioning, the accuracy guarantee itself gets better as L_p gets reduced.*

Destination-Node Aggregate-Frequency Query Estimation

The task of finding the aggregate-frequency of a destination node j is broken down to:

1. Computing the aggregate-frequency of the destination-node j in each partition p , which can be accomplished as described in [3].
2. Computing the sum of all of the computed destination-node aggregate-frequencies together

The first step implies that the time complexity for answering the query is at least n times slower in the worst case than it is without partitioning, where n is the number of partitions after performing the sketch partitioning algorithm. The second step only takes $\mathcal{O}(n)$ operations. Thus, the worst-case time complexity to answer the query is $\mathcal{O}(nhw)$.

Next, we observe how partitioning affects the probabilistic accuracy guarantee.

Theorem 5. *Let $Q_{agg}^-(j)$ be the true aggregate-frequency of destination-node j , and $\overline{Q_{agg}^-(j)}$ be the estimated aggregate-frequency. Let L be the total frequency of edges received so far. Let $\epsilon \in (0, 1)$ be a very small fraction. Then,*

$$P(Q_{agg}^-(j) \leq \overline{Q_{agg}^-(j)} \leq Q_{agg}^-(j) + L \cdot \epsilon) \geq 1 - \left(\frac{n}{h \cdot \epsilon}\right)^w$$

Proof. We note that $P(Q_{agg}^-(j) \leq \overline{Q_{agg}^-(j)}) = 1$ since $\overline{Q_{agg}^-(j)}$ is always an overestimate. We remain to show that

$$P(\overline{Q_{agg}^-(j)} \leq Q_{agg}^-(j) + L \cdot \epsilon) \geq 1 - \left(\frac{n}{h \cdot \epsilon}\right)^w$$

Let $\{p_1, \dots, p_n\}$ be the set of all partitions in the gMatrix and L_i be the total frequency of edges in p_i . In any partition p_i , the probability for a spurious edge (u, v) , $v \neq j$ to be mapped onto any of the rows at the $g_k(j)^{th}$ column, k^{th} layer is $\frac{1}{h}$, so the expected number of spurious edges is $\frac{L_i}{h}$. Let R_k^i be the random variable that represents the number of spurious edges in p_i for the k^{th} hash function. By Markov's inequality,

$$P(R_k^i \geq L_i \epsilon) \leq \frac{E[R_k^i]}{L_i \epsilon} = \frac{1}{h \epsilon} \quad (2.14)$$

Let $R_k = \sum_{i=1}^n R_k^i$ be the total number of spurious edges in all of the n partitions. Combining inequation 2.14 for all $i \in \{1..n\}$,

$$P(R_k \geq L \epsilon) \leq \frac{n}{h \epsilon} \quad (2.15)$$

Therefore,

$$\begin{aligned}
 P(\overline{Q_{agg}^-}(j)) &\leq Q_{agg}^-(j) + L \cdot \epsilon \\
 &= 1 - \prod_{k=1}^w P(R_k \geq L \cdot \epsilon) \\
 &\geq 1 - \left(\frac{n}{h \cdot \epsilon}\right)^w
 \end{aligned} \tag{2.16}$$

□

Remarks. *The probability of the guarantee gets lower as the number of partitions increases, and the guarantee itself never gets any better with partitioning. In other words, partitioning will never improve the accuracy of gMatrix for answering destination-node aggregate-frequency queries.*

Chapter 3

Experiments

We compare the accuracy of gMatrix without partitioning and with partitioning for answering edge frequency estimation queries.

3.1 Implementation

The code used for all of the experiments are implemented in C++. All experiments were performed on Intel Xeon 2GHz 16GB server.

We note that some of the variables mentioned in the earlier chapter are made to be constants for the purpose of the project. These include:

- The data sample used for partitioning is the same in all experiments, i.e. they are only randomly sampled once using Reservoir Sampling [4]. All of the samples' size are 5% of the corresponding dataset size. They are not regenerated on different runtimes of the experiments.
- Statistical variance threshold for filtering source nodes to be used for the partitioning algorithm is set to 100 for all experiments.
- For outlier ratio estimation, data sample is split by 90:10 ratio for all experiments.
- For the first terminating condition of the sketch partitioning algorithm, the minimum number of rows w_0 for any partition is set to 400 in all experiments.
- For the second terminating condition of the sketch partitioning algorithm, the upper-bound constant C is set to 0.1 in all experiments.

3.2 Datasets

The datasets used for the experiments are the same graph streams used for evaluating gMatrix in [3], such as:

- IP-Trace Network Stream [3]
- Twitter Communication Stream [3]
- Friendster Stream with Zipf Frequency distribution [3]

All datasets used for all of the experiments consist of distinct edges only.

3.3 Metrics

3.3.1 Edge Frequency & Node Aggregate-Frequency Estimation

Observed Error [3]

$$observed\ error = \frac{\sum_{i=1}^{|Q|} |\tilde{f}(q_i) - f(q_i)|}{\sum_{i=1}^{|Q|} f(q_i)}$$

where $Q = \{q_i\}$ is the set of queries, \tilde{f} is the estimated frequency, and f is the actual frequency.

Average Relative Error [1]

$$average\ relative\ error = \frac{1}{|Q|} \sum_{i=1}^{|Q|} \frac{|\tilde{f}(q_i) - f(q_i)|}{f(q_i)}$$

where $Q = \{q_i\}$ is the set of queries, \tilde{f} is the estimated frequency, and f is the actual frequency.

Effective Queries [1]

Both Observed Error and Average Relative Error may be biased if the frequencies of the queries vary a lot, so we define queries as "effective" if the relative error is not exceeding T . The percentage of effective queries is computed by,

$$effective\ queries = \frac{|\{q | \frac{|\tilde{f}(q) - f(q)|}{f(q)} \leq T, q \in Q\}|}{|Q|} \cdot 100\%$$

where $Q = \{q_i\}$ is the set of queries, \tilde{f} is the estimated frequency, and f is the actual frequency.

In the project, T is a constant and is set to 5, as also is the case in [1].

3.3.2 Heavy Hitter Edge Estimation

False Positive Rate [3]

The metric measures the percentage of edges that are misclassified as heavy-hitter edges.

$$false\ positive\ rate(\%) = \frac{false\ heavyhitter\ edges}{true\ heavyhitter\ edges} \cdot 100\%$$

3.4 Results

3.4.1 Edge Frequency Estimation

Friendster dataset Partitioning reduces the observed error and the average relative error of the estimations, as seen on figures 3.1-3.3. The lower errors are possible because partitioning reduces the number of queries with high relative errors (queries q such that $\frac{\tilde{f}(q)-f(q)}{f(q)} > T$), as the number of effective queries themselves are actually lower with partitioning as seen on figure 3.5.

Twitter dataset Partitioning does not seem to be effective at all as seen on figures 3.7-3.12.

IP-Trace dataset Partitioning reduces the observed error and the average relative error of the estimations, as seen on figures 3.13-3.15. Unlike the results on the Friendster dataset, the number of effective queries after partitioning are only significantly lower on $h = 1000$ and is actually higher than without partitioning on $h = 4000$.

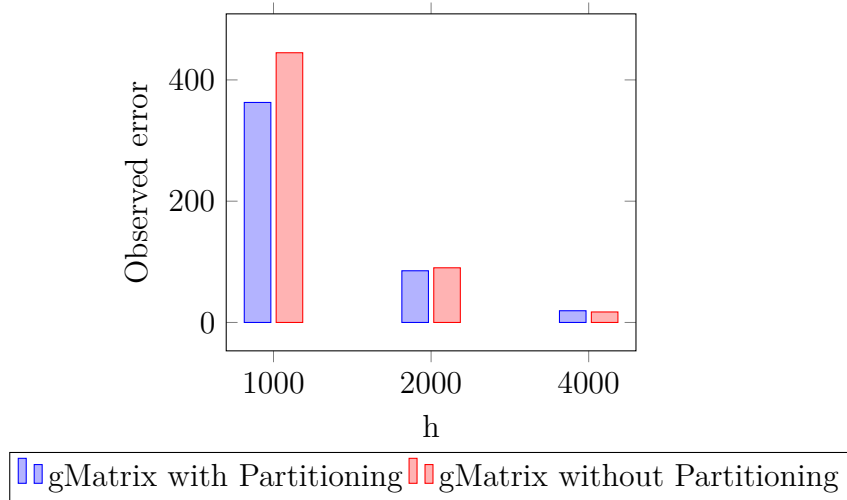


Fig. 3.1 Observed error on Friendster dataset considered over 1 million random edges

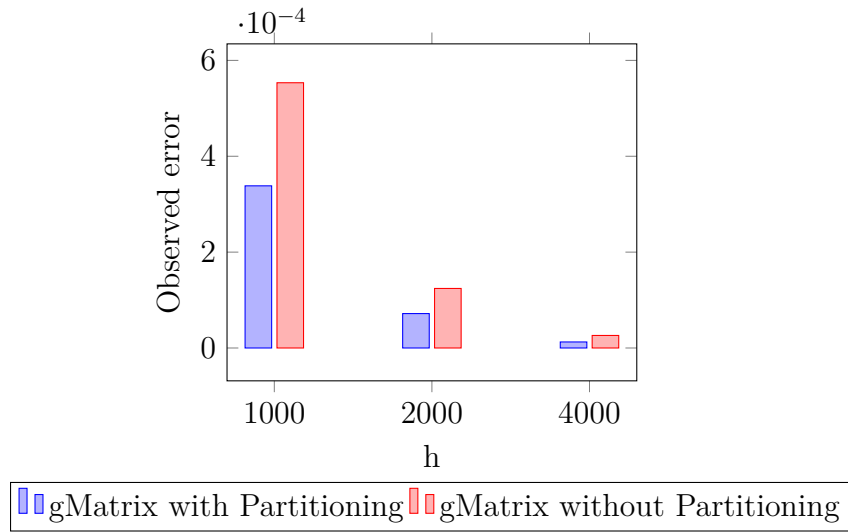


Fig. 3.2 Observed error on Friendster dataset considered over top-500 highest frequency edges

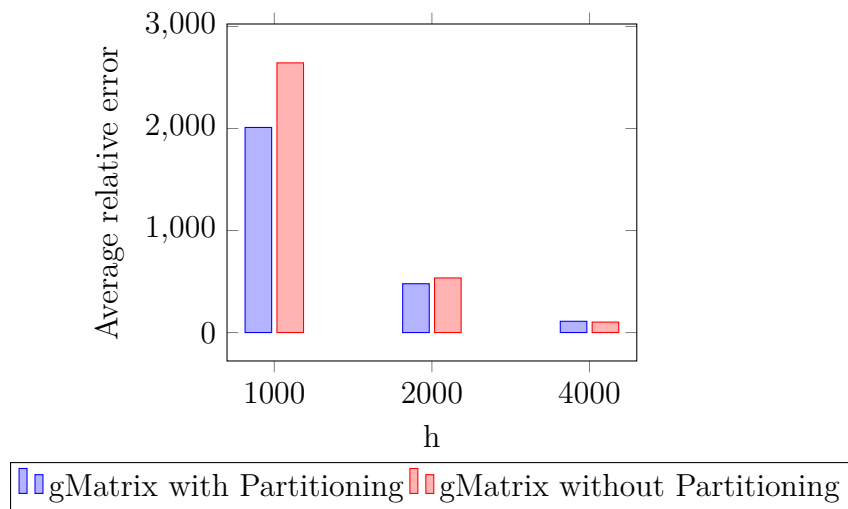


Fig. 3.3 Average relative error on Friendster dataset considered over 1 million random edges

Experiments

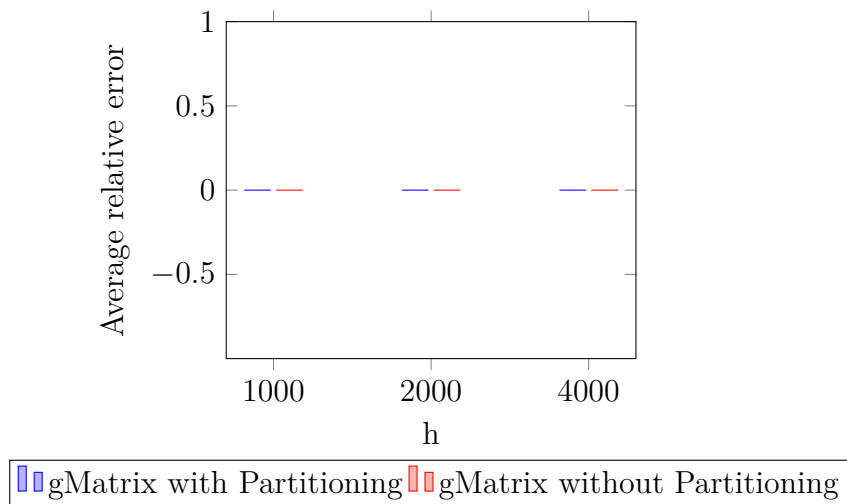


Fig. 3.4 Average relative error on Friendster dataset considered over top-500 highest frequency edges

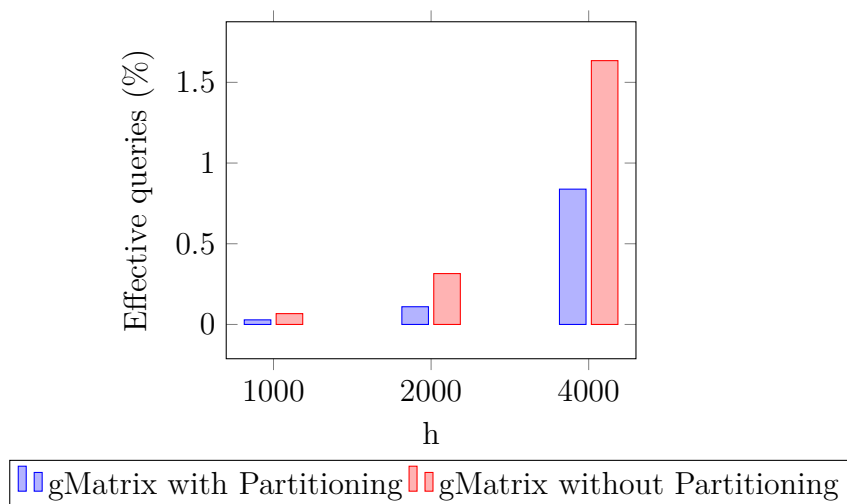


Fig. 3.5 Effective queries percentage on Friendster dataset considered over 1 million random edges

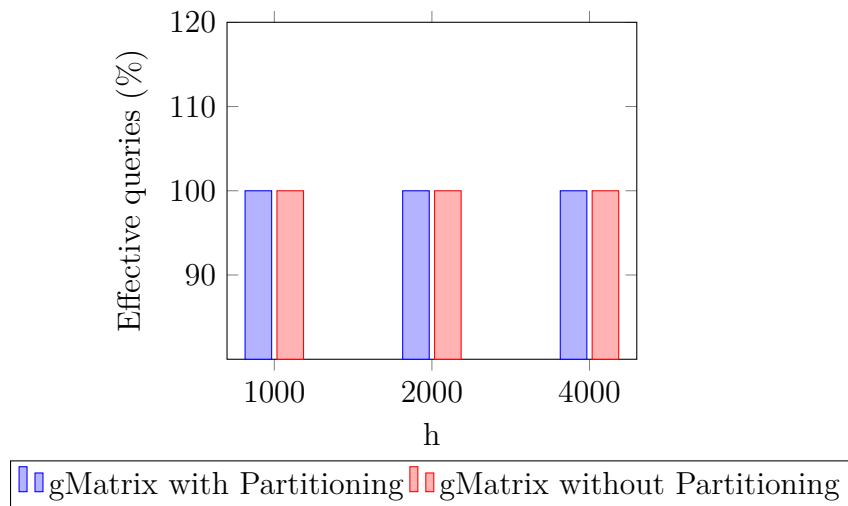


Fig. 3.6 Effective queries percentage on Friendster dataset considered over top-500 highest frequency edges

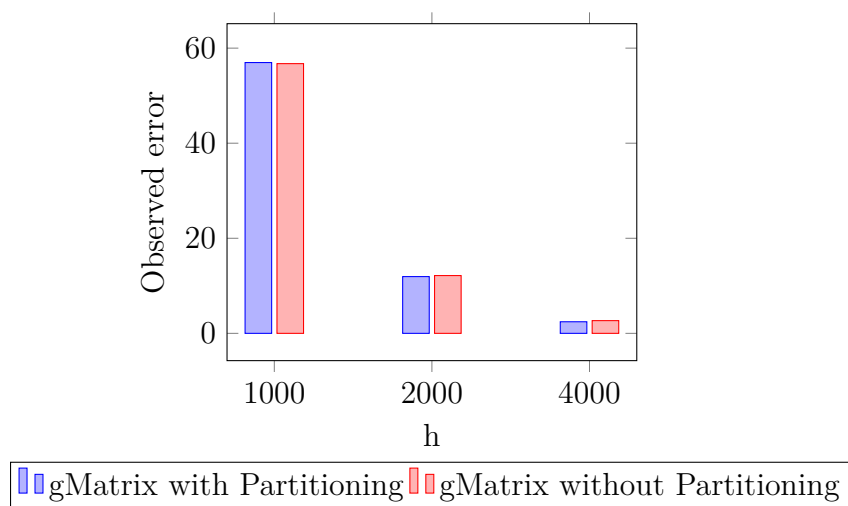


Fig. 3.7 Observed error of Twitter dataset considered over 1 million random edges

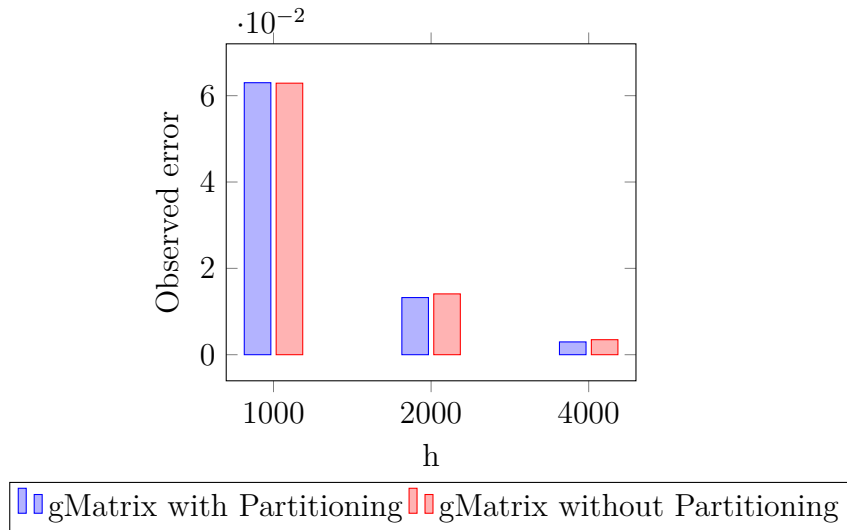


Fig. 3.8 Observed error of Twitter dataset considered over top-500 highest frequency edges

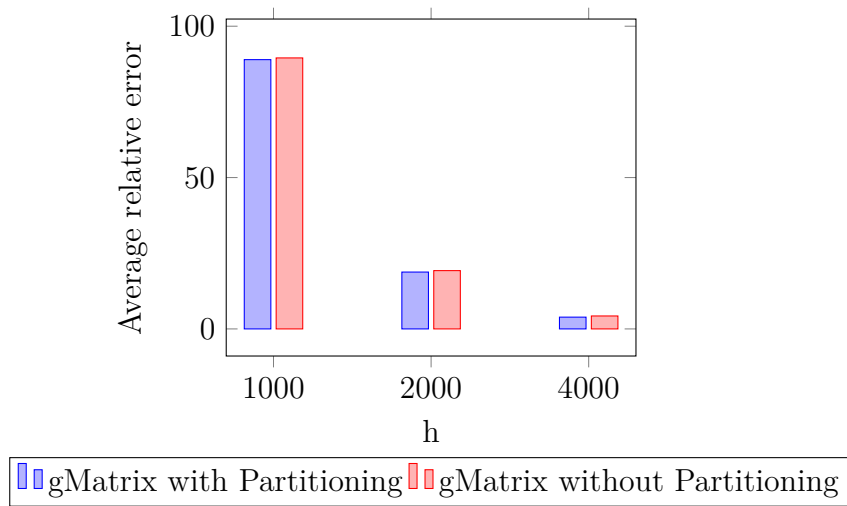


Fig. 3.9 Average relative error on Twitter dataset considered over 1 million random edges

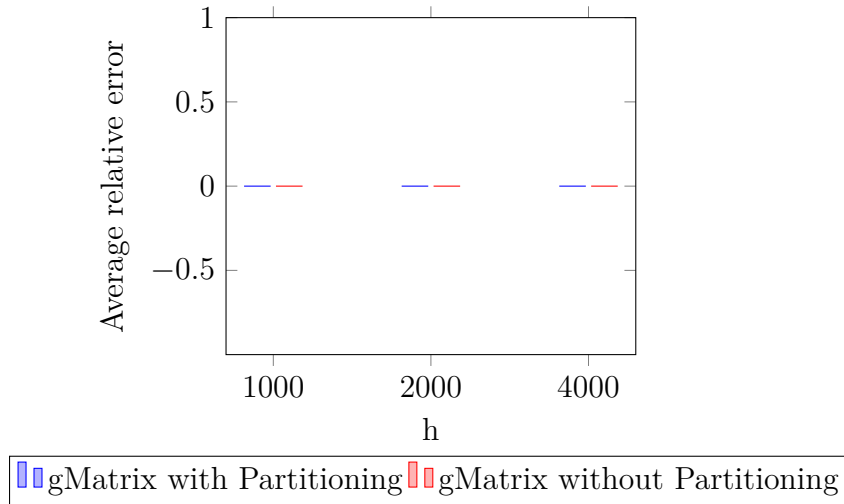


Fig. 3.10 Average relative error on Twitter dataset considered over top-500 highest frequency edges

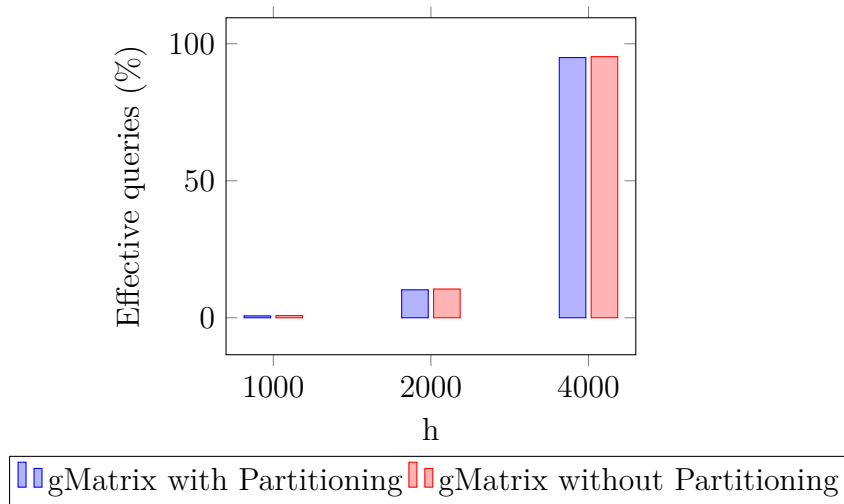


Fig. 3.11 Effective queries percentage on Twitter dataset considered over 1 million random edges

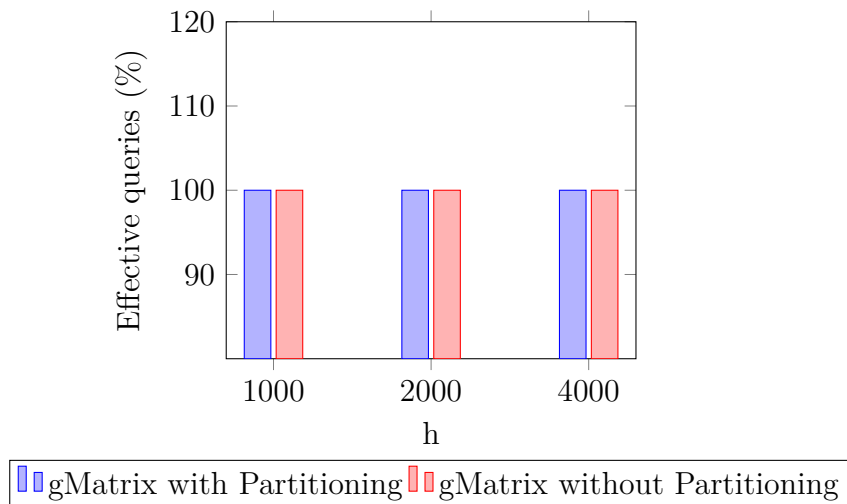


Fig. 3.12 Effective queries percentage on Twitter dataset considered over top-500 highest frequency edges

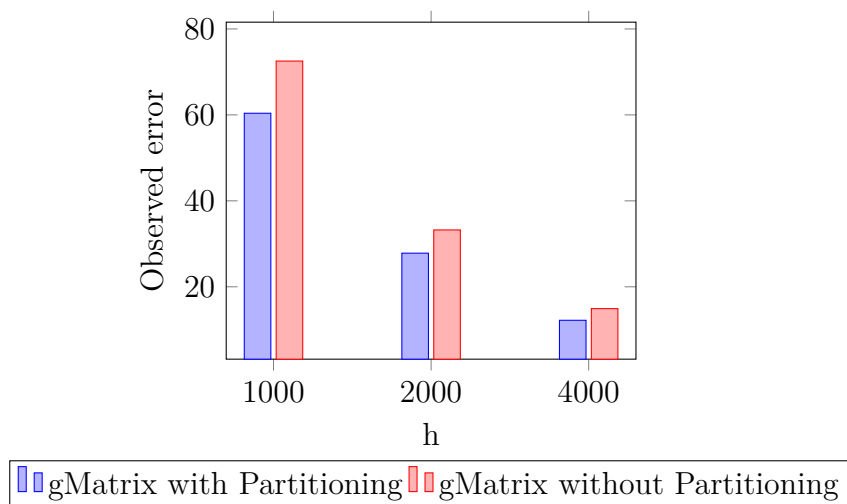


Fig. 3.13 Observed error of IP-Trace dataset considered over 1 million random edges

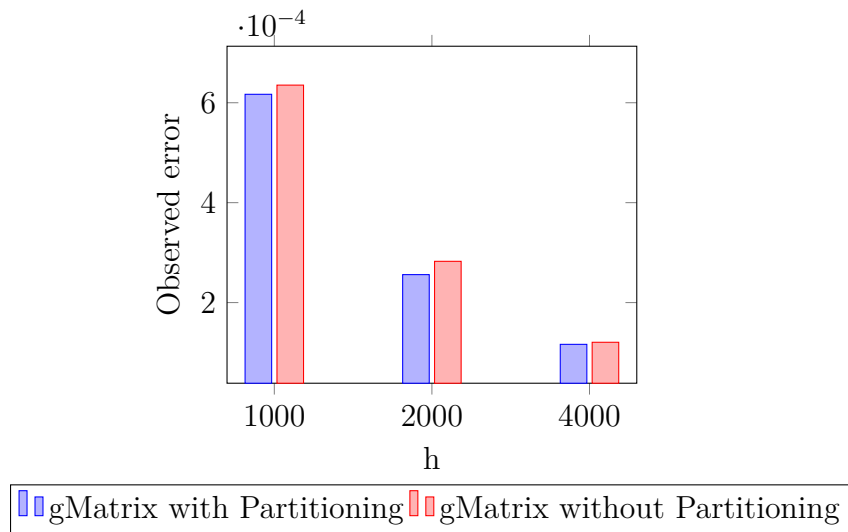


Fig. 3.14 Observed error of IP-Trace dataset considered over top-500 highest frequency edges

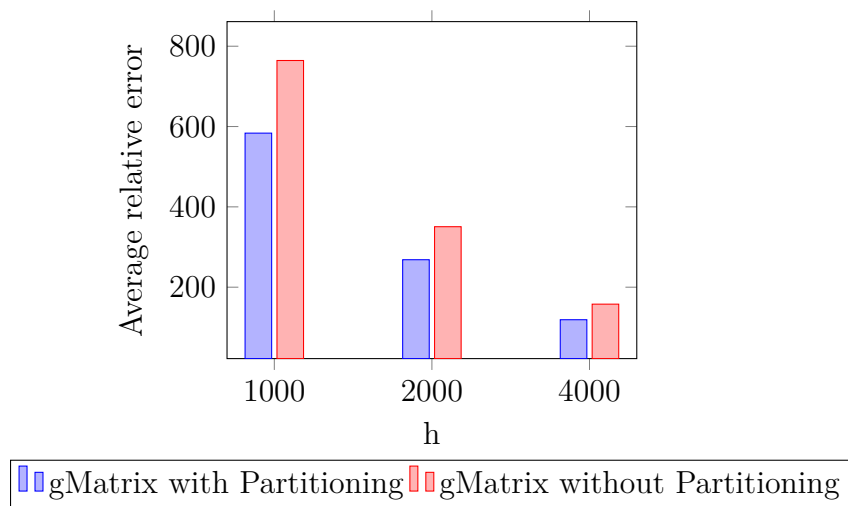


Fig. 3.15 Average relative error on IP-Trace dataset considered over 1 million random edges

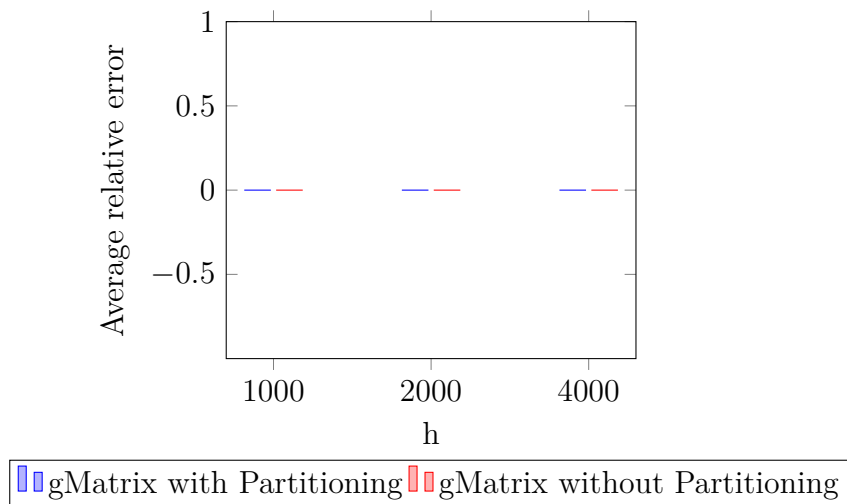


Fig. 3.16 Average relative error on IP-Trace dataset considered over top-500 highest frequency edges

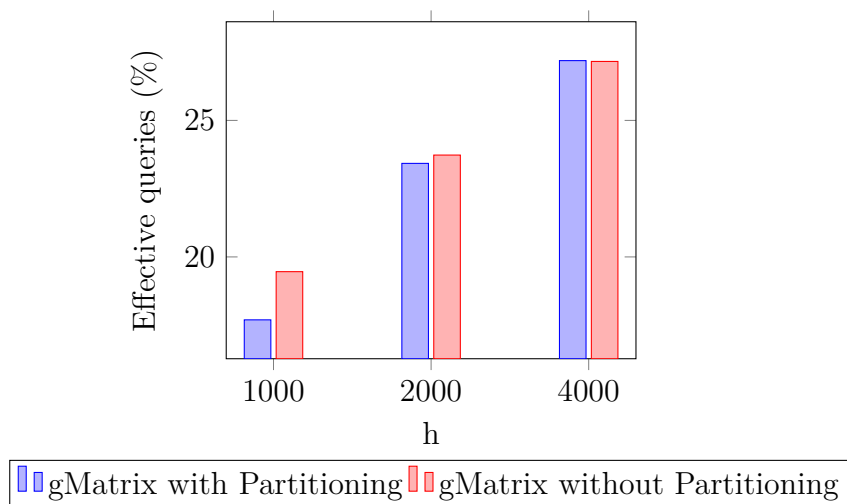


Fig. 3.17 Effective queries percentage on IP-Trace dataset considered over 1 million random edges

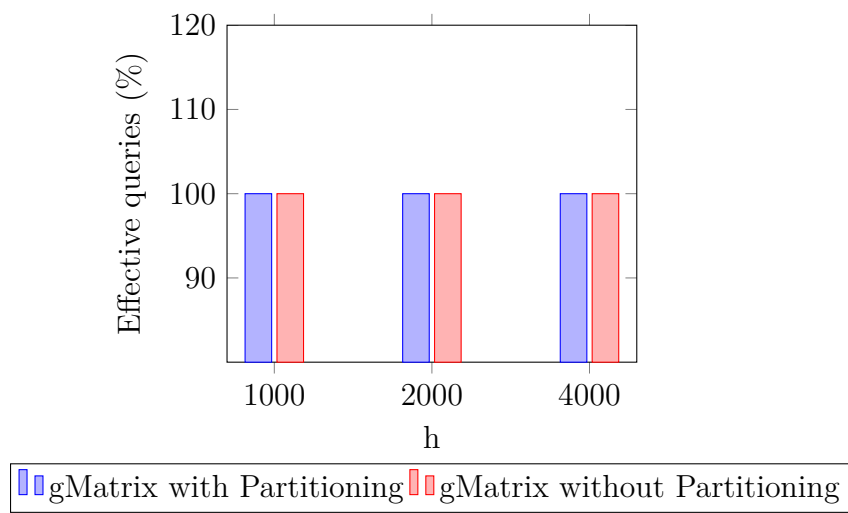


Fig. 3.18 Effective queries percentage on IP-Trace dataset considered over top-500 highest frequency edges

3.4.2 Heavy Hitter Edge Estimation

Figures 3.19 and 3.20 show that after partitioning, the sketch is more sensitive to low frequency thresholds as the false positive rates are much higher at frequency threshold of 0.01% of the total stream frequency. However, it performs as good at frequency threshold percentages of 0.1% and 1%, with no false positives produced at all.

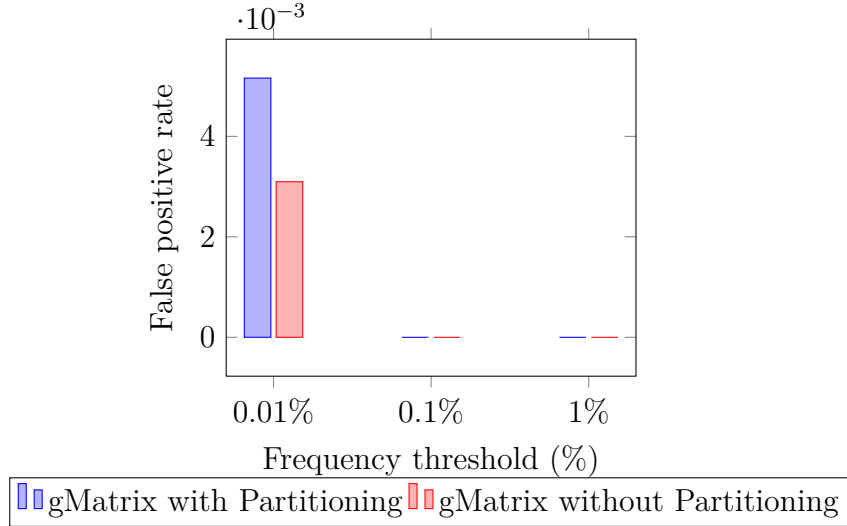


Fig. 3.19 False positive rate of heavy hitter edge estimation on IP-Trace dataset with respect to varying frequency threshold percentages of the total stream frequency

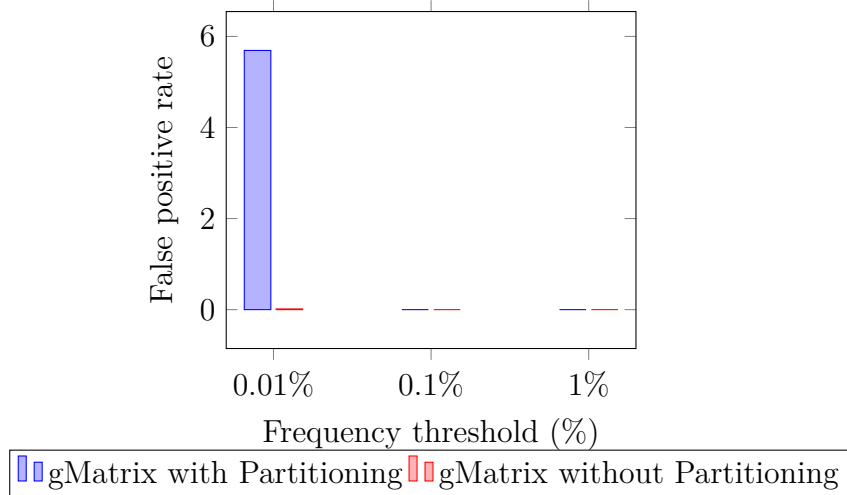


Fig. 3.20 False positive rate of heavy hitter edge estimation on Friendster dataset with respect to varying frequency threshold percentages of the total stream frequency

3.4.3 Node Aggregate-Frequency Estimation

On the Twitter and IP-Trace dataset, the observed error of source-node aggregate-frequency estimations is lower after partitioning as seen on figure 3.23 and figure 3.25. However, on the Friendster dataset, figure 3.21 shows that it is higher after partitioning, which is unexpected. According to Theorem 4, an ideal partitioning scheme improves accuracy guarantee, although reduces the probability of the guarantee. We speculate that the higher observed error on the Friendster dataset is due to the high-skewness of the dataset, which makes the partitioning scheme less effective.

Figures 3.22, 3.24 and 3.26 show the observed error is higher after partitioning in destination-node aggregate-frequency estimations, which is expected since Theorem 5 shows that partitioning does not improve the accuracy guarantee itself and only worsens the probability of the guarantee.

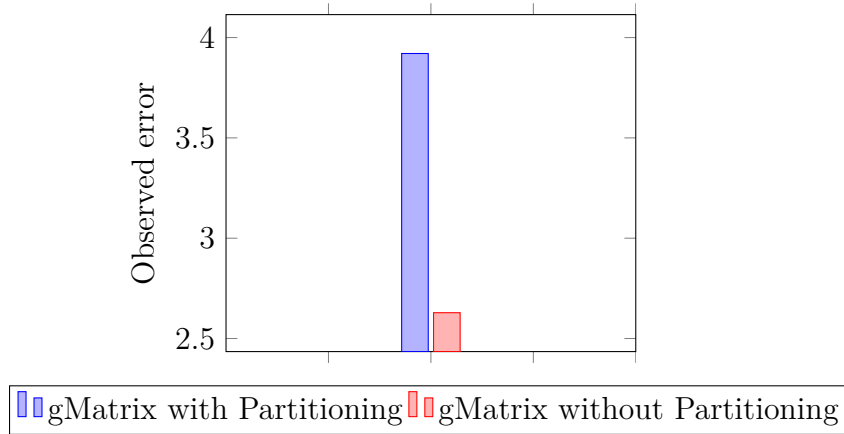


Fig. 3.21 Observed error of source-node aggregate-frequency queries on top-500 node aggregate-frequencies of the Friendster dataset

Experiments

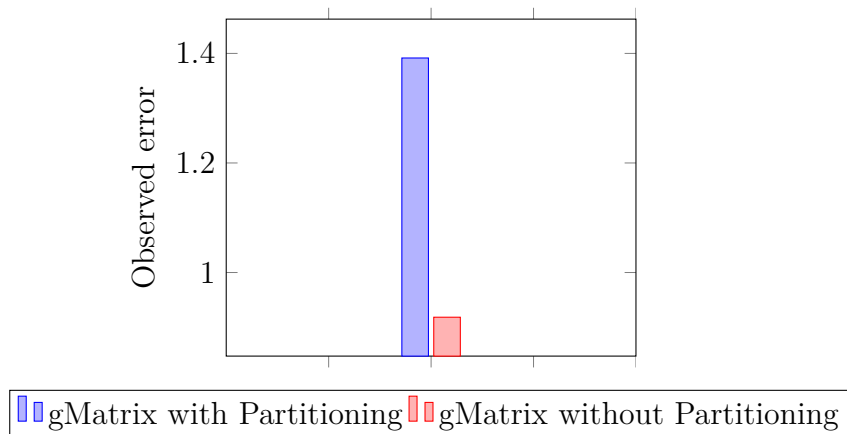


Fig. 3.22 Observed error of destination-node aggregate-frequency queries on top-500 node aggregate-frequencies of the Friendster dataset

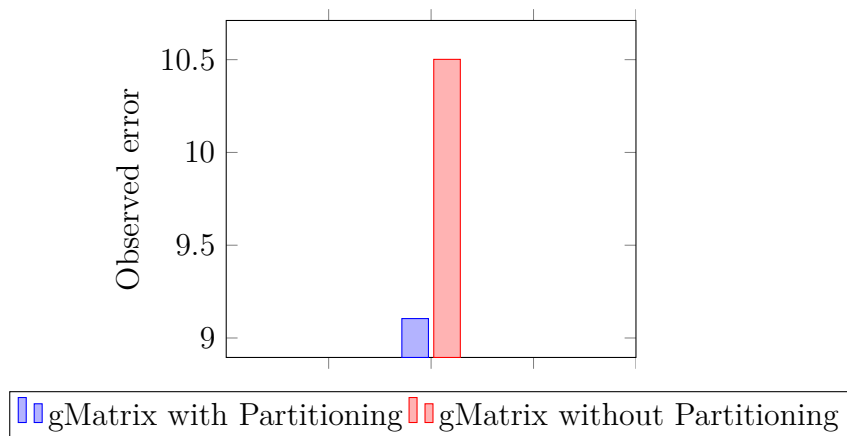


Fig. 3.23 Observed error of source-node aggregate-frequency queries on top-500 node aggregate-frequencies of the Twitter dataset

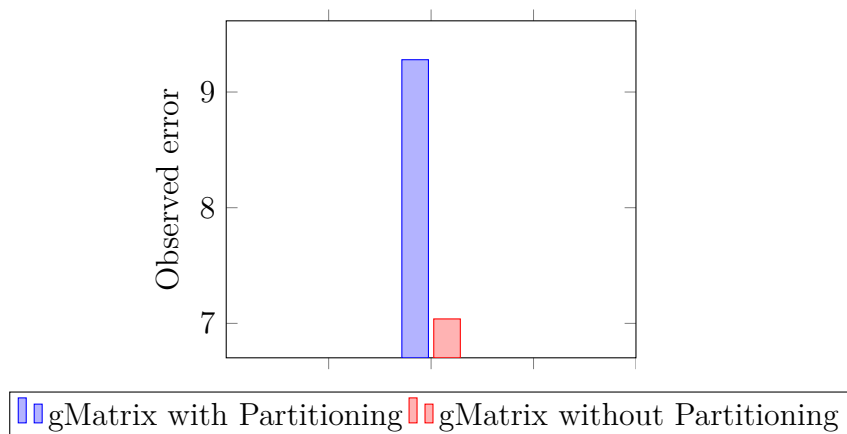


Fig. 3.24 Observed error of destination-node aggregate-frequency queries on top-500 node aggregate-frequencies of the Twitter dataset

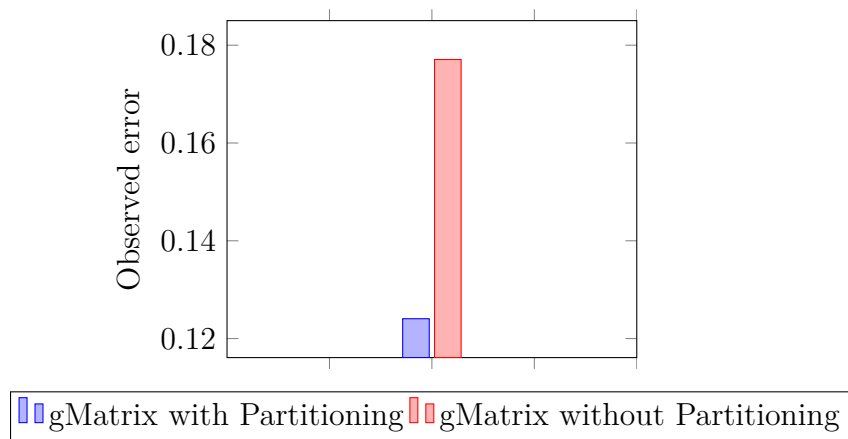


Fig. 3.25 Observed error of source-node aggregate-frequency queries on top-500 node aggregate-frequencies of the IP-Trace dataset

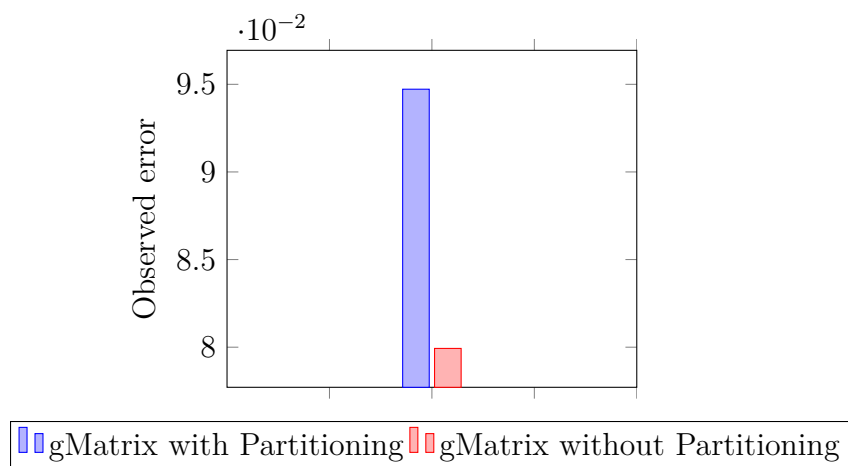


Fig. 3.26 Observed error of destination-node aggregate-frequency queries on top-500 node aggregate-frequencies of the IP-Trace dataset

Chapter 4

Conclusions

4.1 Conclusion

To summarize all contributions made so far, we have proposed optimizations to the partitioning method, analyzes how the method impacts how the gMatrix sketch answers various query types and how the performance and accuracy guarantees change, and conducted experiments on some datasets.

The results show that the gSketch partitioning method generally improves gMatrix accuracy on query types such as edge frequency and source-node aggregate frequency estimations. In addition, such queries are performed with the same time complexity too.

However, the gSketch partitioning method causes gMatrix to be less accurate and to be slower on some query types such as the destination-node aggregate frequency estimation and the heavy-hitter edge queries.

Finally, we note that whether the partitioning method is beneficial for use depends on the intended use-case.

4.2 Possible Improvements

Experiments on how the partitioning method affects gMatrix on other query types supported by the gMatrix are yet to be conducted. In addition, experiments on the effects of varying the partitioning parameters on the gMatrix are yet to be conducted too. Finally, experiments on other datasets may produce new insights.

References

- [1] P. Zhao, C. C. Aggarwal, and M. Wang, “gsketch: On query estimation in graph streams,” *CoRR*, vol. abs/1111.7167, 2011.
- [2] G. Cormode and S. Muthukrishnan, “An improved data stream summary: the count-min sketch and its applications,” *Journal of Algorithms*, vol. 55, no. 1, pp. 58–75, 2005.
- [3] A. Khan and C. Aggarwal, “Toward query-friendly compression of rapid graph streams,” *Social Network Analysis and Mining*, vol. 7, no. 1, p. 23, 2017.
- [4] J. S. Vitter, “Random sampling with a reservoir,” *ACM Transactions on Mathematical Software (TOMS)*, vol. 11, no. 1, pp. 37–57, 1985.