

# A Partitioning Scheme for the gMatrix Graph Stream Sketch



**Eric Leonardo Lim**

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) of the*  
*Nanyang Technological University*

AY 2017/2018



## **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 is where you write your abstract ...



# Table of contents

<b>List of figures</b>	<b>ix</b>
<b>1 Introduction</b>	<b>1</b>
1.1 Background . . . . .	1
1.2 Objectives . . . . .	1
1.3 Relevant Literatures . . . . .	1
<b>2 Design and Analysis</b>	<b>3</b>
2.1 Preliminaries . . . . .	3
2.1.1 CountMin Sketch . . . . .	3
2.1.2 gSketch . . . . .	3
2.1.3 gMatrix . . . . .	3
2.2 Partitioning Scheme . . . . .	3
2.3 Edge Frequency Queries . . . . .	5
2.4 Heavy-hitter Edge Queries . . . . .	7
2.5 Node Aggregate-Frequency Queries . . . . .	9
2.5.1 Source-Node Aggregate-Frequency Queries . . . . .	9
2.5.2 Destination-Node Aggregate-Frequency Queries . . . . .	10
<b>3 Experiments</b>	<b>13</b>
3.1 System Description . . . . .	13
3.2 Datasets . . . . .	13
3.3 Metrics . . . . .	13
3.3.1 Average Relative Error [3] . . . . .	14
3.3.2 Effective Queries [3] . . . . .	14
3.4 Results . . . . .	15
3.4.1 Edge Frequency Estimation . . . . .	15
3.4.2 Heavy Hitter Edge Estimation . . . . .	25

4	Conclusion	27
---	------------	----

	References	29
--	------------	----



# List of figures

3.1	Observed error on Friendster dataset considered over 1 million random edges . . . . .	15
3.2	Observed error on Friendster dataset considered over top-500 highest frequency edges . . . . .	15
3.3	Average relative error on Friendster dataset considered over 1 million random edges . . . . .	16
3.4	Average relative error on Friendster dataset considered over top-500 highest frequency edges . . . . .	16
3.5	Effective queries percentage on Friendster dataset considered over 1 million random edges . . . . .	17
3.6	Effective queries percentage on Friendster dataset considered over top-500 highest frequency edges . . . . .	17
3.7	Observed error of Twitter dataset considered over 1 million random edges	19
3.8	Observed error of Twitter dataset considered over top-500 highest frequency edges . . . . .	19
3.9	Average relative error on Twitter dataset considered over 1 million random edges . . . . .	20
3.10	Average relative error on Twitter dataset considered over top-500 highest frequency edges . . . . .	20
3.11	Effective queries percentage on Twitter dataset considered over 1 million random edges . . . . .	21
3.12	Effective queries percentage on Twitter dataset considered over top-500 highest frequency edges . . . . .	21
3.13	Observed error of IP-Trace dataset considered over 1 million random edges	22
3.14	Observed error of IP-Trace dataset considered over top-500 highest frequency edges . . . . .	22

3.15	Average relative error on IP-Trace dataset considered over 1 million random edges . . . . .	23
3.16	Average relative error on IP-Trace dataset considered over top-500 highest frequency edges . . . . .	23
3.17	Effective queries percentage on IP-Trace dataset considered over 1 million random edges . . . . .	24
3.18	Effective queries percentage on IP-Trace dataset considered over top-500 highest frequency edges . . . . .	24
3.19	False positive rate of heavy hitter edge estimation on IP-Trace dataset	25
3.20	False positive rate of heavy hitter edge estimation on Friendster dataset	25

# 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 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.

### 1.2 Objectives

The objective of this work is to improve gMatrix's [2] accuracy using the idea of partitioning given data sample [3]. Experiments are to be carried on multiple dat

### 1.3 Relevant Literatures

Previous work [3] has introduced a way to utilize the Count-min sketch [1], which is commonly used to approximate the frequency of events in a general data stream, to approximate the frequency of edges in a graph stream. Furthermore, it [3] has introduced the gSketch, which is a partitioning scheme that exploits typical local structural properties within real world graph datasets by building multiple localized Count-min sketches [1] to summarize a graph stream. It has been shown that the

partitioning scheme is able to improve overall accuracy when tested on several large graph datasets.

A recent work [2] has introduced the gMatrix, which is a sketch that generalizes the Count-min sketch for 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 [3].

Applying a partitioning scheme identical to the gSketch [3] into the gMatrix [2] may improve its overall accuracy. However, there are no previous works that have shown if it really is the case and the extent of the improvement.

# Chapter 2

## Design and Analysis

### 2.1 Preliminaries

#### 2.1.1 CountMin Sketch

#### 2.1.2 gSketch

#### 2.1.3 gMatrix

### 2.2 Partitioning Scheme

We implement the idea of partitioning from [3] onto the gMatrix.

A data sample is needed to perform the partitioning. In the experiment, we perform reservoir sampling on the original dataset to obtain the data sample. The size of the data sample is 5% of the original size of the dataset.

The gMatrix is divided into partitions. Each partition 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 one of the partitions if and only if  $u$  is in the set of source vertices associated with the partition. The idea of the partitioning is to group edges with similar frequencies in the same partition to improve sketch accuracy.

The statistical variances of the frequencies of edges (from the data sample) with common source vertex are computed. We select source vertices which have variances not exceeding a certain threshold (variance of 100 is used in the experiment). An outlier partition is reserved to store frequencies of edges in the dataset whose source vertex is either not found in the data sample or are not suitable for the purpose of the

partitioning, i.e. the statistical variance of the frequencies of all edges with the same source vertex in the data sample is exceeding the threshold.

Determining the right outlier partition size is very critical for the effectiveness of the partitioning, even more so than in partitioning the gSketch[3]. We determine the outlier partition size ratio by estimating the outlier ratio from the data sample. We do so by first splitting the data sample into two; the first consists of 90% of the whole data sample, and the other split consists of the rest; then, after selecting the suitable source vertices from the first split, the ratio of outliers in the second split is computed, which is essentially the ratio of the number of source vertices in the split not found among the previously-selected source vertices, to the size of the split. Then, it is assigned as the outlier partition size ratio.

The partitioning algorithm is no different than in [3].

---

**Algorithm 1** Sketch-Partitioning (Data Sample:  $D$ )

---

```

1: Create a root node  $S$  of the partitioning tree as an active node;
2:  $S.sides \leftarrow h$ ;
3:  $S.depth \leftarrow d$ ;
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 overall
   relative error;
7:    $S_1.rows = S_2.rows = \frac{S.rows}{2}$ ;
8:    $L \leftarrow L \setminus \{S\}$ ;
9:   for  $i \in [1..2]$  do
10:    if ( $S_i.sides \geq w_0$ ) and ( $\sum_{v \in V_S} \tilde{d}(v) \leq C \cdot S_i.sides$ ) then
11:       $L \leftarrow L \cup S_i$ ;
12:    else
13:      Construct the localized sketch  $S_i$ ;
```

---

The suitable source vertices are sorted based on non-decreasing average frequency of edges emanating from them  $f_v/d_v$  in the data sample, where  $f_v$  is the sum of frequencies of all edges in the data sample with source vertex  $v$ , and  $d_v$  is the out-degree of vertex  $v$ . The sorted source vertex set is then recursively split into two, and each split minimizes the overall relative error  $E$ .

$$E = \sum_{v \in S_1} \frac{d_v \cdot F_{S_1}}{f_v/d_v} + \sum_{v \in S_2} \frac{d_v \cdot F_{S_1}}{f_v/d_v}$$

The recursion stops when either the size of the partition becomes small enough, i.e. the number of rows is less than a user-specified threshold  $r_0$  (which is fixed at 100

in this experiment), or when the probability of collision in any particular cell in the sketch can be bounded above by a user-specified threshold  $C$ , which is the case when the density of distinct edges,  $\sum_{v \in s} d_v / (S.rows \cdot S.cols)$ , is also bounded above by  $C$ , as proven in [3].

## 2.3 Edge Frequency Queries

### Description

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 the minimum of  $V_P(g_k(i), g_k(j), k)$  for all  $k \in \{1..w\}$ , where  $V_P$  is the value of the cell in the partition  $p$ .

### Analysis

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 is 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 1.** *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,

$$P(X_k \geq L_p \cdot \epsilon) \leq \frac{E[X_k]}{L_p \cdot \epsilon} = \frac{2^d}{h^2 \epsilon} \quad (2.1)$$

The second possible case of spurious edges  $(u, v)$  are those for which  $u \neq i$  but  $v = j$ . The probability that 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,

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

The third possible case of spurious edges  $(u, v)$  are those for which  $u = i$  but  $v \neq j$ . The probability that 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,

$$P(Z_k \geq h B_p(i) \cdot \epsilon) \leq \frac{E[Z_k]}{h B_p(i) \cdot \epsilon} = \frac{1}{h^2 \epsilon} \quad (2.3)$$

Combining inequations 2.1, 2.2, 2.3,

$$P(X_k + Y_k + Z_k \geq L_p \cdot \epsilon + A_p(j) \cdot h \cdot \epsilon + B_p(i) \cdot h \cdot \epsilon) \leq \frac{2^{d+1} + 1}{h^2 \cdot \epsilon} \quad (2.4)$$



Therefore,

$$\begin{aligned}
P(\overline{Q(i, j)}) &\leq Q(i, j) + L_p \cdot \epsilon + A_p(i) \cdot h \cdot \epsilon + B_p(i) \cdot h \cdot \epsilon \\
&= 1 - \prod_{k=1}^w P(X_k + Y_k + Z_k \geq L_p \cdot \epsilon + A_p(j) \cdot h \cdot \epsilon + B_p(i) \cdot h \cdot \epsilon) \\
&\geq 1 - \left(\frac{2^{d+1} + 1}{h^2 \cdot \epsilon}\right)^w
\end{aligned} \tag{2.5}$$

□

**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.*

## 2.4 Heavy-hitter Edge Queries

### Description

Obtaining edges with frequencies at least  $F$  in the graph stream is accomplished by first obtaining a set of edges with frequencies at least  $F$  from each partition in the gMatrix and then taking the union of each of them.

### Analysis

**Theorem 2.** *Let  $(i, j)$  be an edge and  $p$  be the partition associated with  $i$ . 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  $\frac{1}{h \cdot \frac{h}{2^d}} = \frac{2^d}{h^2}$ . 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,

$$P(X_k \geq \frac{\overline{Q(i, j)} - F}{3}) \leq \frac{E[X_k]}{\frac{\overline{Q(i, j)} - F}{3}} = \frac{3 \cdot 2^d L_p}{h^2 (\overline{Q(i, j)} - F)} \quad (2.6)$$

The second possible case of spurious edges  $(u, v)$  are those for which  $u \neq i$  but  $v = j$ . The probability that 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,

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

The third possible case of spurious edges  $(u, v)$  are those for which  $u = i$  but  $v \neq j$ . The probability that 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,

$$P(Z_k \geq \frac{\overline{Q(i, j)} - F}{3}) \leq \frac{E[Z_k]}{\frac{\overline{Q(i, j)} - F}{3}} = \frac{3 \cdot B_p(i)}{h (\overline{Q(i, j)} - F)} \quad (2.8)$$

Combining inequations 2.6, 2.7, 2.8,

$$P(X_k + Y_k + Z_k \geq \overline{Q(i, j)} - F) \leq \frac{3 \cdot 2^d (L_p + h \cdot A_p(j)) + 3 \cdot h \cdot B_p(i)}{h^2 \cdot (\overline{Q(i, j)} - F)} \quad (2.9)$$

Therefore,

$$\begin{aligned} & P(Q(i, j) \geq F) \\ &= 1 - \prod_{k=1}^w P(X_k + Y_k + Z_k \geq \overline{Q(i, j)} - F) \\ &\geq 1 - \left( \frac{3 \cdot 2^d (L_p + h \cdot A_p(j)) + 3 \cdot h \cdot B_p(i)}{h^2 \cdot (\overline{Q(i, j)} - F)} \right)^w \end{aligned} \quad (2.10)$$

□

## 2.5 Node Aggregate-Frequency Queries

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 computed differently.

### 2.5.1 Source-Node Aggregate-Frequency Queries

#### Description

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

1. Finding the associated partition  $p$  of  $i$
2. Finding the aggregate frequency of the source node  $i$  at the partition  $p$

#### Analysis

**Theorem 3.** *Let  $(i, j)$  be an edge and  $p$  be the partition associated with  $i$ . Let  $Q_{agg}^+(i)$  be the true aggregate-frequency of source-node  $i$ , and  $\overline{Q_{agg}^+(i)}$  be its estimated aggregate-frequency. Let  $p$  be a partition associated with  $i$  and  $L_p$  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  $(g_k(i), \cdot, k)$  is  $\frac{1}{h} = \frac{2^d}{h}$ , so the expected number of spurious edges that are mapped into  $(g_k(i), \cdot, k)$  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,

$$P(R_k \geq L_p \epsilon) \leq \frac{E[R_k]}{L_p \epsilon} = \frac{2^d}{h \epsilon} \quad (2.11)$$

Therefore,

$$\begin{aligned}
P(\overline{Q_{agg}^+}(i)) &\leq Q_{agg}^+(i) + L_p \cdot \epsilon \\
&= 1 - \prod_{k=1}^w P(R_k \geq L_p \cdot \epsilon) \\
&\geq 1 - \left(\frac{2^d}{h \cdot \epsilon}\right)^w
\end{aligned} \tag{2.12}$$

□

**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.

## 2.5.2 Destination-Node Aggregate-Frequency Queries

### Description

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  $p$
2. Computing the sum of all of the computed destination-node aggregate-frequencies together

### Analysis

**Theorem 4.** 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  $(\cdot, g_k(j), k)$  is  $\frac{1}{h}$ , so the expected number of spurious edges that

are mapped into  $(\cdot, g_k(j), k)$  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,

$$P(R_k^i \geq L_i \epsilon) \leq \frac{E[R_k^i]}{L_i \epsilon} = \frac{1}{h \epsilon} \quad (2.13)$$

Let  $R_k = \sum_{i=1}^n R_k^i$ . Combining inequation 2.13 for all  $i \in \{1..n\}$ ,

$$P(R_k \geq L \epsilon) \leq \frac{n}{h \epsilon} \quad (2.14)$$

Therefore,

$$\begin{aligned} P(\overline{Q_{agg}^-}(j) \leq Q_{agg}^-(j) + L \cdot \epsilon) \\ &= 1 - \prod_{k=1}^w P(R_k \geq L \cdot \epsilon) \\ &\geq 1 - \left(\frac{n}{h \cdot \epsilon}\right)^w \end{aligned} \quad (2.15)$$

□

**Remarks.** *The probability of the guarantee gets lower as the number of partitions increases, and the guarantee itself never gets any better with partitioning.*



# Chapter 3

## Experiments

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

### 3.1 System Description

The code is implemented in C++ and experiments were performed on Intel Xeon 2GHz 16GB server.

### 3.2 Datasets

The datasets used to evaluate the sketches are graph streams consisting of distinct edges only. The datasets are:

- IP-Trace Network Stream [2]
- Twitter Communication Stream [2]
- Friendster Stream with Zipf Frequency distribution [2]

### 3.3 Metrics

Observed Error [2]

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

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

### 3.3.1 Average Relative Error [3]

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

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

### 3.3.2 Effective Queries [3]

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 this experiment,  $T$  is set to 5, which is the same as [3].



## 3.4 Results

### 3.4.1 Edge Frequency Estimation

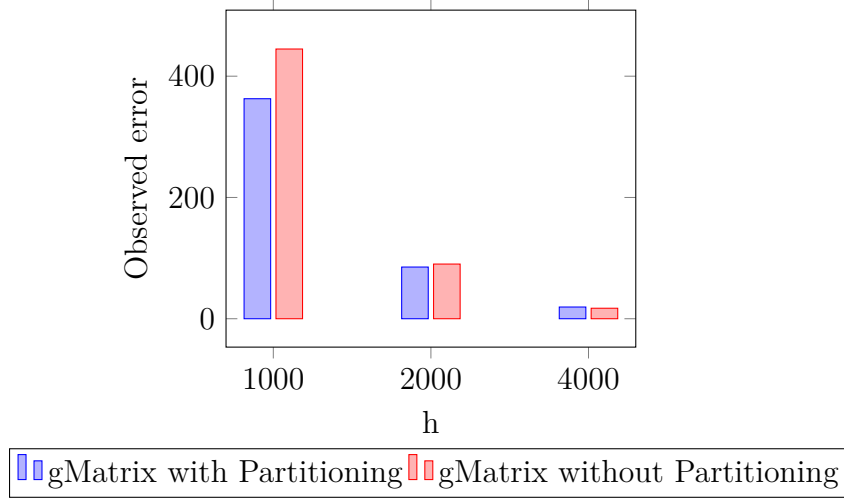


Fig. 3.1 Observed error on Friendster dataset considered over 1 million random edges

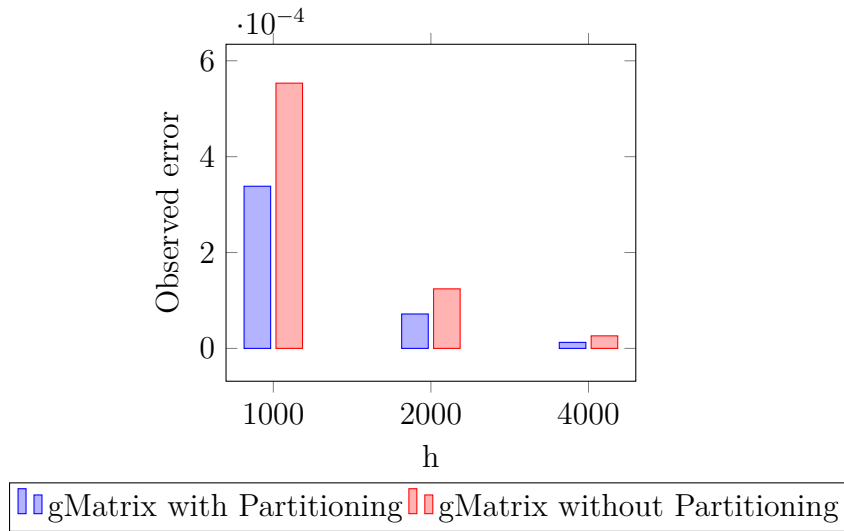


Fig. 3.2 Observed error on Friendster dataset considered over top-500 highest frequency edges

**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

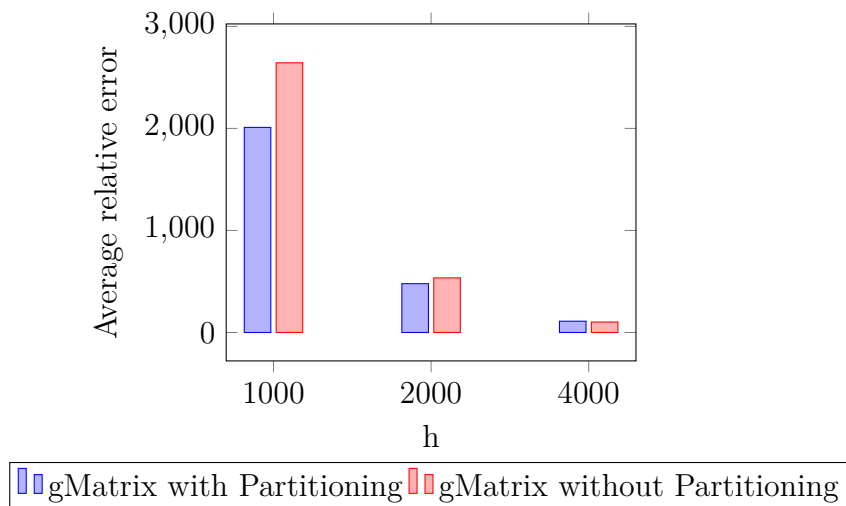


Fig. 3.3 Average relative error on Friendster dataset considered over 1 million random edges

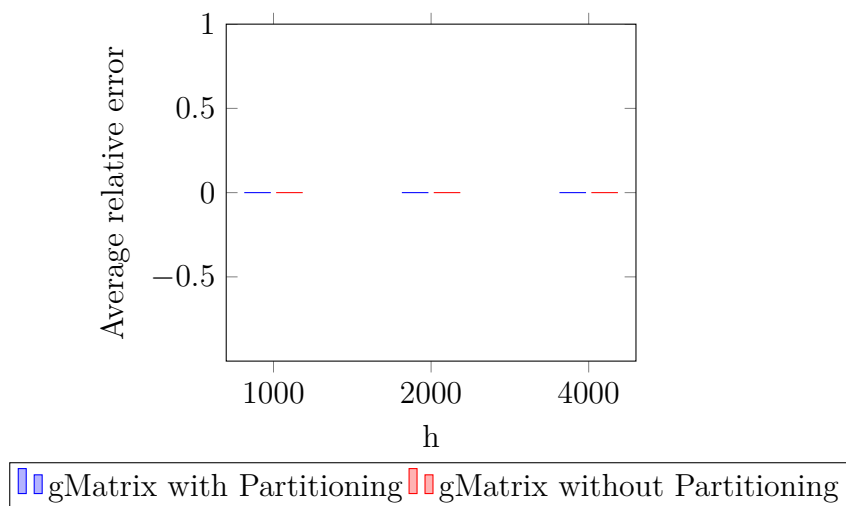


Fig. 3.4 Average relative error on Friendster dataset considered over top-500 highest frequency edges

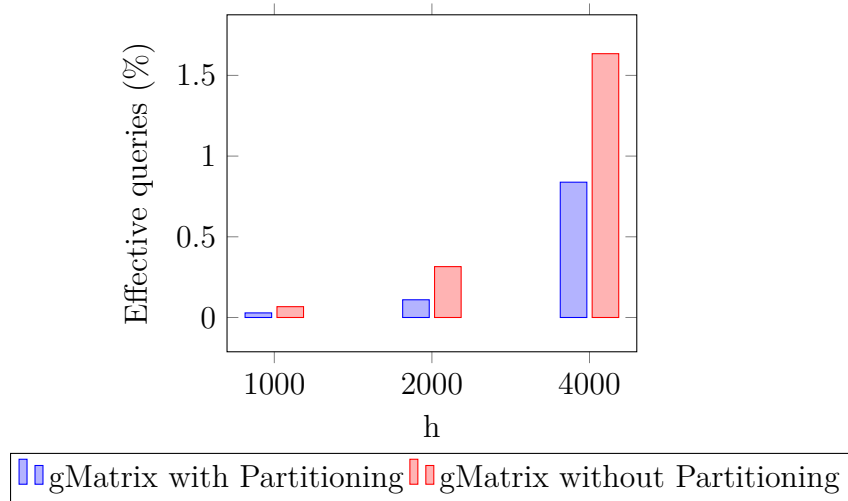


Fig. 3.5 Effective queries percentage on Friendster dataset considered over 1 million random edges

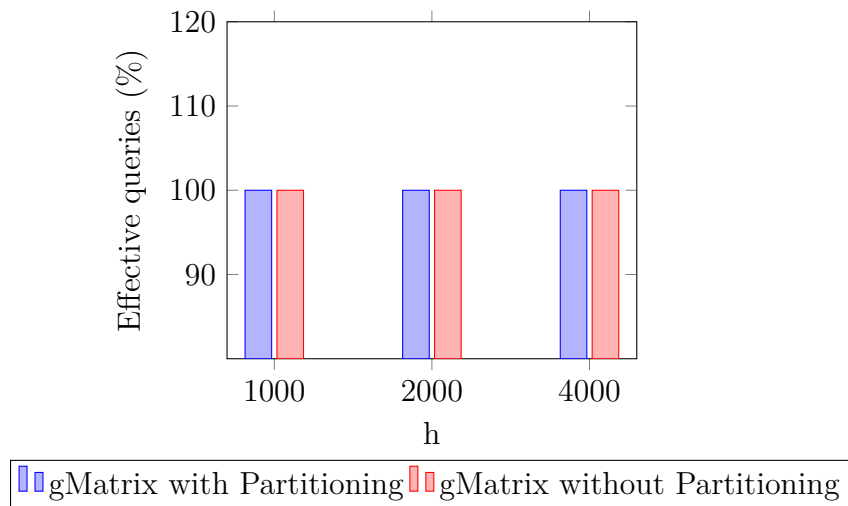


Fig. 3.6 Effective queries percentage on Friendster dataset considered over top-500 highest frequency edges

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.

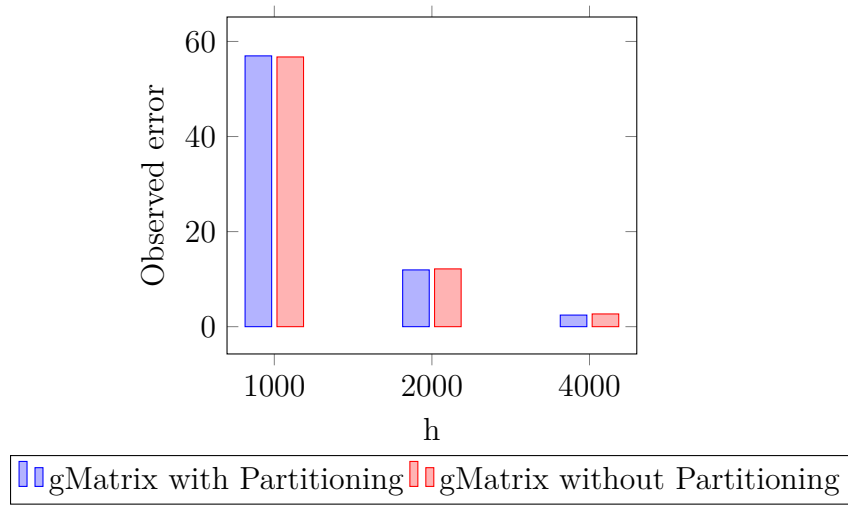


Fig. 3.7 Observed error of Twitter dataset considered over 1 million random edges

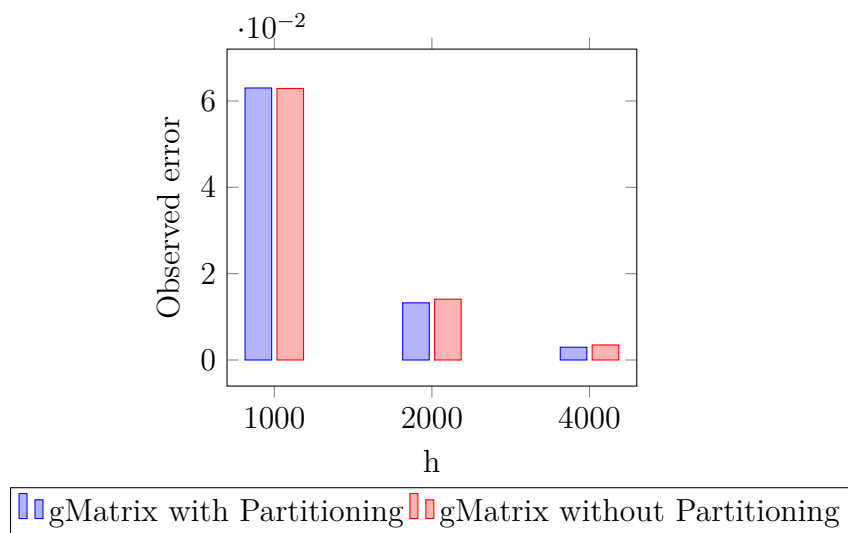


Fig. 3.8 Observed error of Twitter dataset considered over top-500 highest frequency edges

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

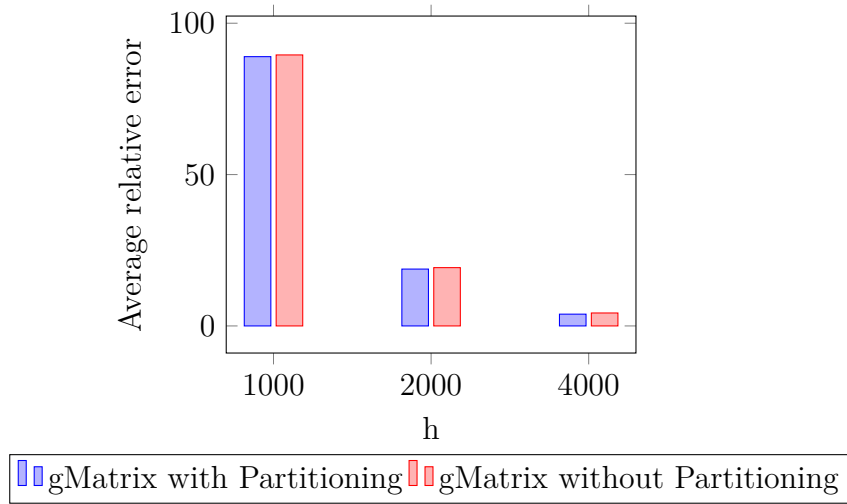


Fig. 3.9 Average relative error on Twitter dataset considered over 1 million random edges

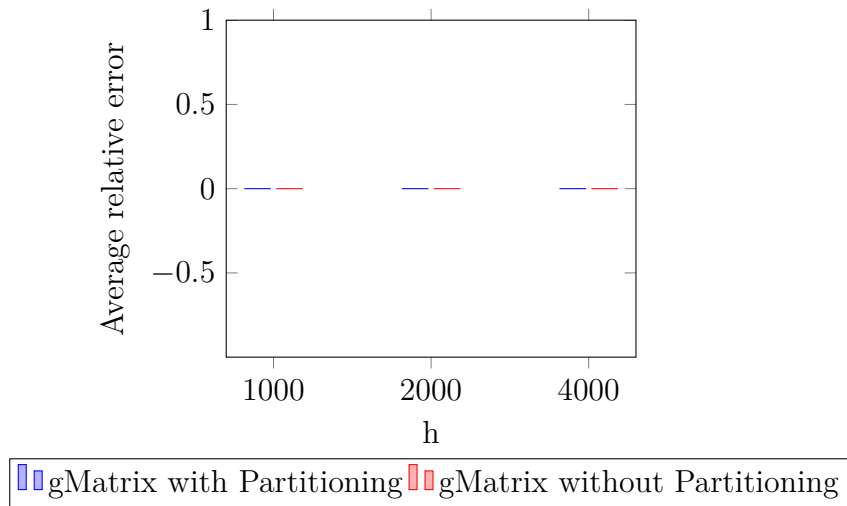


Fig. 3.10 Average relative error on Twitter dataset considered over top-500 highest frequency edges

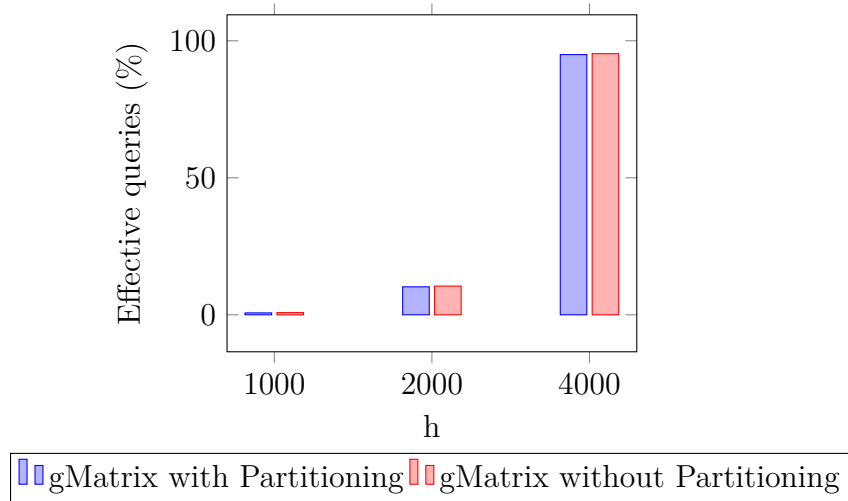


Fig. 3.11 Effective queries percentage on Twitter dataset considered over 1 million random edges

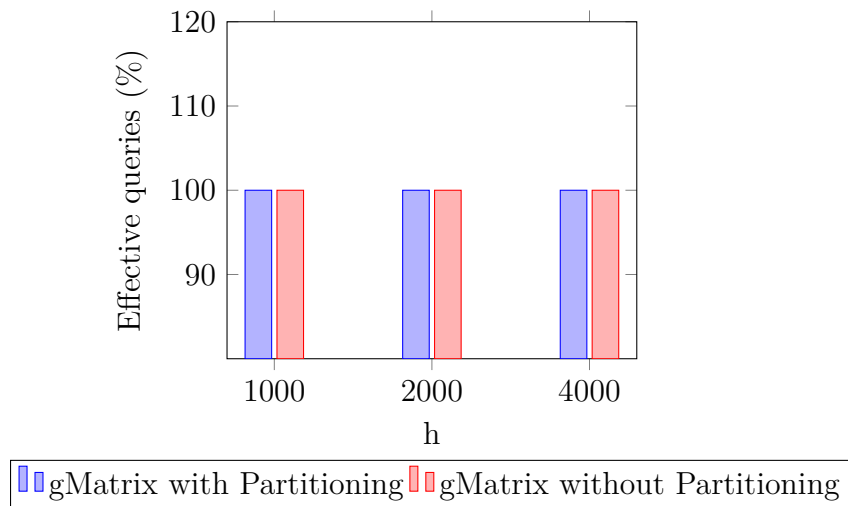


Fig. 3.12 Effective queries percentage on Twitter dataset considered over top-500 highest frequency edges

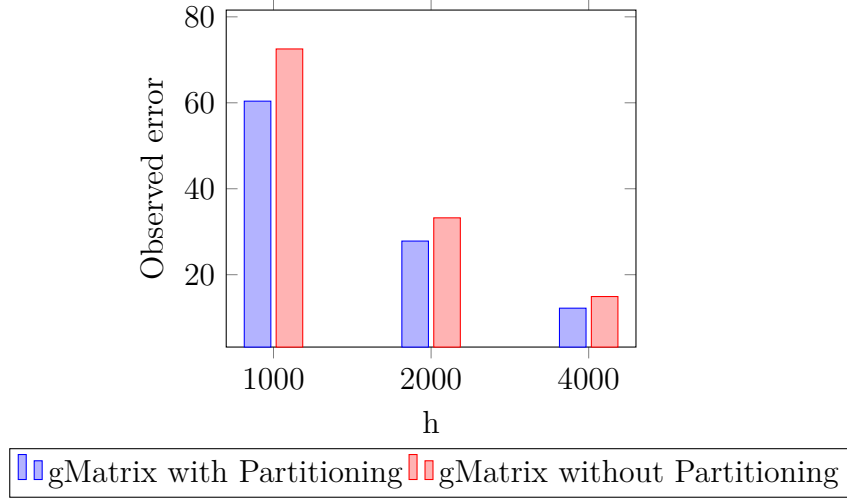


Fig. 3.13 Observed error of IP-Trace dataset considered over 1 million random edges

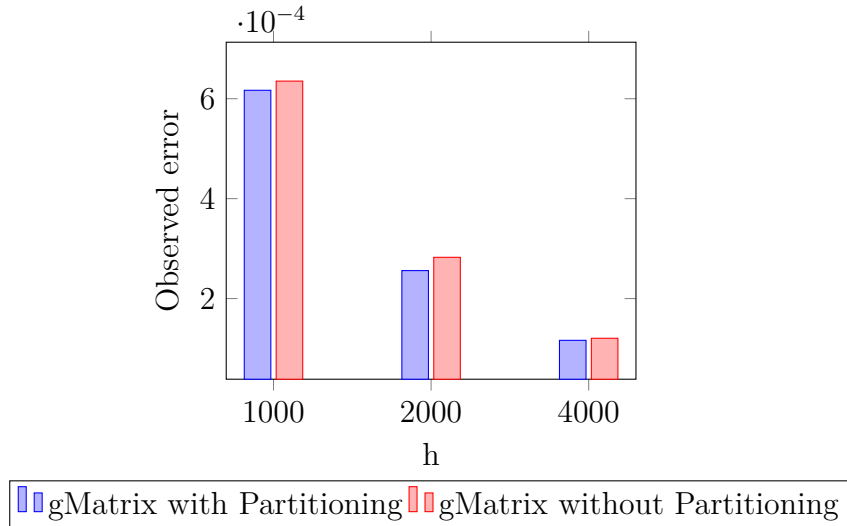


Fig. 3.14 Observed error of IP-Trace dataset considered over top-500 highest frequency edges

**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  $h = 4000$ .



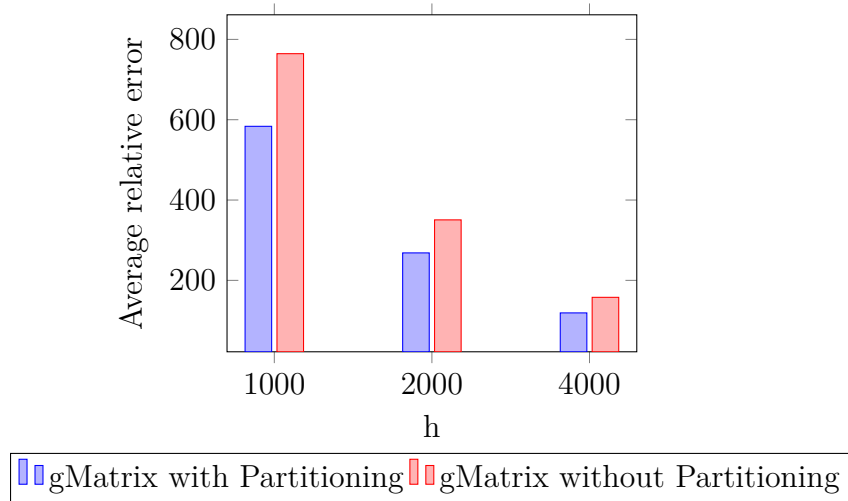


Fig. 3.15 Average relative error on IP-Trace dataset considered over 1 million random edges

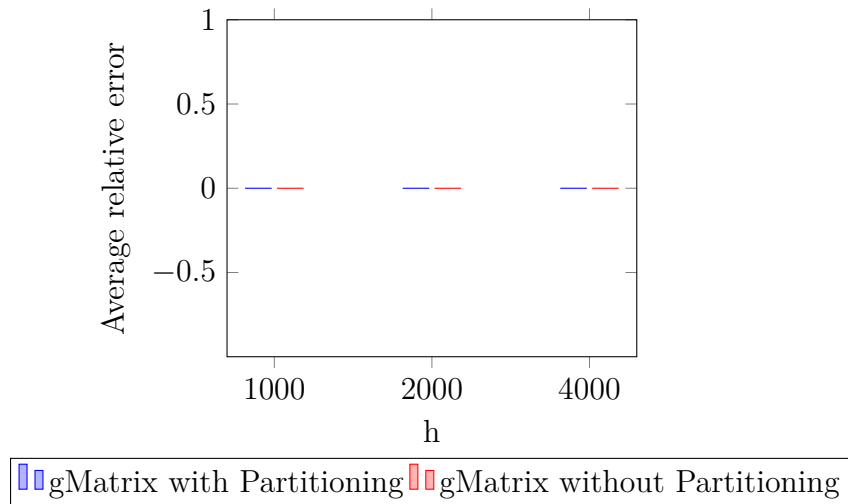


Fig. 3.16 Average relative error on IP-Trace dataset considered over top-500 highest frequency edges

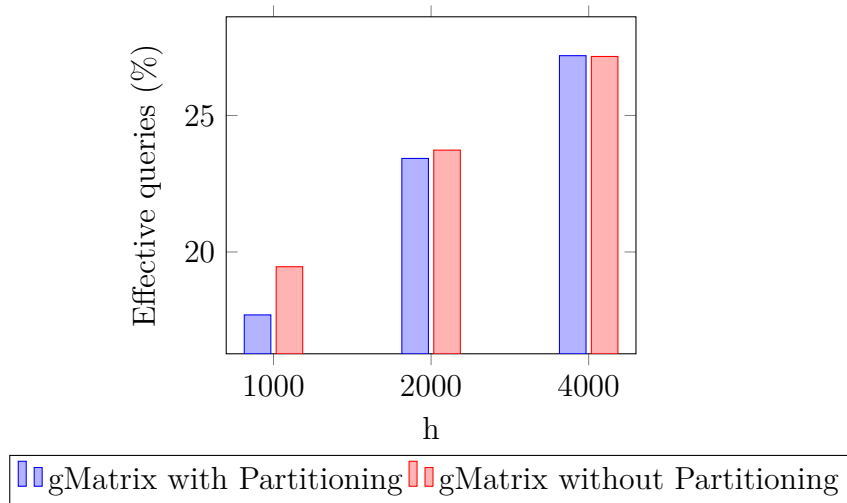


Fig. 3.17 Effective queries percentage on IP-Trace dataset considered over 1 million random edges

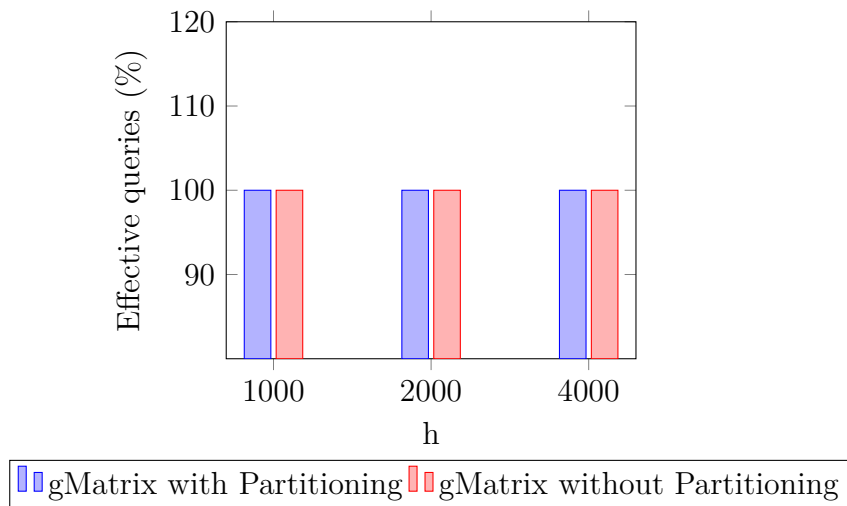


Fig. 3.18 Effective queries percentage on IP-Trace dataset considered over top-500 highest frequency edges

### 3.4.2 Heavy Hitter Edge Estimation

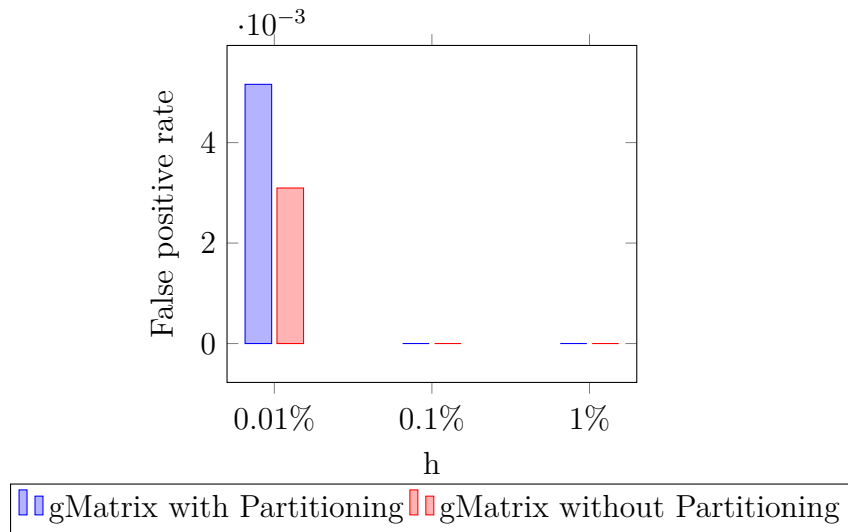


Fig. 3.19 False positive rate of heavy hitter edge estimation on IP-Trace dataset

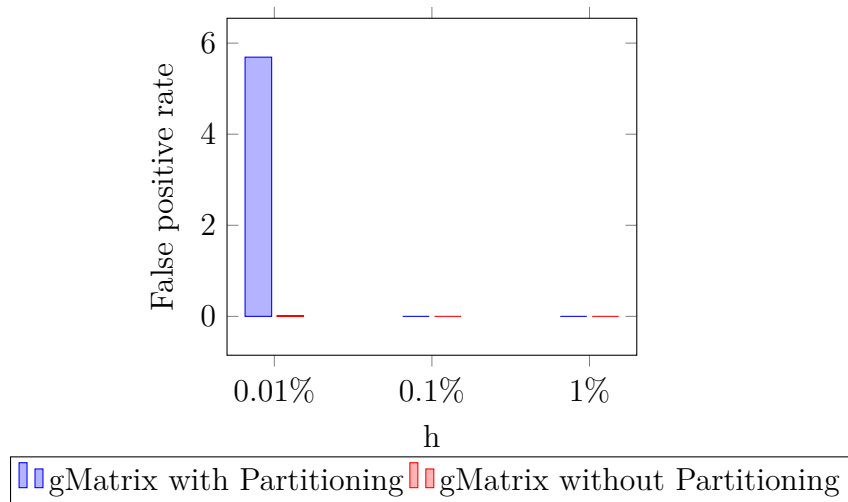


Fig. 3.20 False positive rate of heavy hitter edge estimation on Friendster dataset



## Chapter 4

## Conclusion



# References

- [1] Cormode, G. and Muthukrishnan, S. (2005). An improved data stream summary: the count-min sketch and its applications. *Journal of Algorithms*, 55(1):58–75.
- [2] Khan, A. and Aggarwal, C. (2017). Toward query-friendly compression of rapid graph streams. *Social Network Analysis and Mining*, 7(1):23.
- [3] Zhao, P., Aggarwal, C. C., and Wang, M. (2011). gsketch: On query estimation in graph streams. *CoRR*, abs/1111.7167.