



PER-FLOW CARDINALITY ESTIMATION BASED ON VIRTUAL LOGLOG  
SKETCHING

BY

ZEYU ZHOU

THESIS

Submitted in partial fulfillment of the requirements  
for the degree of Master of Science in Electrical and Computer Engineering  
in the Graduate College of the  
University of Illinois at Urbana-Champaign, 2016

Urbana, Illinois

Adviser:

Professor Bruce Hajek

# ABSTRACT

Flow cardinality estimation is the problem of estimating the number of distinct elements in a data flow, often with a stringent memory constraint. It has wide applications in network traffic measurement and in database systems. The virtual LogLog algorithm proposed recently by Xiao, Chen, Chen and Ling estimates the cardinalities of a large number of flows with a compact memory. The purpose of this thesis is to explore two new perspectives on the estimation process of this algorithm. Firstly, we propose and investigate a family of estimators that generalizes the original vHLL estimator and evaluate the performance of the vHLL estimator compared to other estimators in this family. Secondly, we propose an alternative solution to the estimation problem by deriving a maximum-likelihood estimator. Empirical evidence from both perspectives suggests the near-optimality of the vHLL estimator for per-flow estimation, analogous to the near-optimality of the HLL estimator for single-flow estimation.

# ACKNOWLEDGMENTS

I thank Prof. Bruce Hajek, for introducing this topic to me and guiding me through the completion of this thesis with great patience. I thank my parents, for their unconditional love and support.

# CONTENTS

Chapter 1	INTRODUCTION . . . . .	1
1.1	Flow cardinality estimation . . . . .	1
1.2	Motivation and contributions . . . . .	3
1.3	Thesis overview . . . . .	4
Chapter 2	BACKGROUND AND PRELIMINARIES . . . . .	5
2.1	Single-flow cardinality estimation . . . . .	5
2.2	Per-flow cardinality estimation with memory sharing . . . . .	8
2.3	A per-flow estimator performance metric – Weighted square error . . . . .	12
Chapter 3	VIRTUAL LOGLOG ESTIMATOR WITH PARAMETER $\theta$ . . . . .	16
3.1	The $LL_\theta$ estimator for single-flow estimation . . . . .	16
3.2	The $vLL_\theta$ estimator for per-flow estimation . . . . .	20
3.3	Experimental performance evaluation . . . . .	25
Chapter 4	MAXIMUM LIKELIHOOD ESTIMATOR . . . . .	26
4.1	Formulation . . . . .	26
4.2	Derivation of the likelihood function . . . . .	27
4.3	Implementation . . . . .	32
4.4	Aside: MLE for single-flow estimation . . . . .	36
4.5	Experimental performance evaluation . . . . .	37
Chapter 5	CONCLUSIONS AND FUTURE WORK . . . . .	40
	REFERENCES . . . . .	42

# Chapter 1

## INTRODUCTION

### 1.1 Flow cardinality estimation

Today’s Internet is flooded with data. The measuring and monitoring of Internet data traffic has led to many useful applications but it is also challenging due to the sheer volume and speed of the traffic. *Flow cardinality estimation* is one of the fundamental problems in network traffic measurement and it is the problem addressed by this thesis.

In general, we define *cardinality estimation* to be the problem of estimating the number of *distinct* elements  $n$  of a given multiset<sup>1</sup>  $\mathcal{S}$ . For example, if  $\mathcal{S} = \{1, 2, 3, 2, 1\}$ , then  $n = 3$ . The problem is often considered in the stream model — we are only allowed to observe each element in  $\mathcal{S}$  once, then the element is discarded forever.

Flow cardinality estimation, in its essence, is cardinality estimation in the context of network traffic measurement, where we define a *flow* to be a multiset of data packets defined by certain properties observed on a network link over a period of time. The purpose of flow cardinality estimation is to estimate the number of *distinct* packets in the multiset, where the distinction is made based on some properties such as values of certain packet header fields of interest.

For example, we can consider a *per-source flow* that contains all the packets from the same *source address* and use cardinality estimation to estimate the number of *distinct destination addresses* among them. In this case, the particular source address is known as the flow’s identifier (*flow ID*). Some commonly used flow identifiers are source/destination address, source/destination port number, protocol type, or a combination of them (such as source-destination pair).

**Challenges:** The challenges of flow cardinality estimation are mainly twofold: time and space constraints. First of all, high data transfer rates in some communication networks

---

<sup>1</sup>A multiset is a generalization of a set, with the difference that a multiset can have duplicate elements.

today make it infeasible to spend much time on processing each element. According to [1] and [2], on a typical OC-768<sup>2</sup> backbone network link with 40 Gbps traffic speed, the time available to process each packet is at best about tens of nanoseconds, corresponding to no more than a hundred elementary operations.

Secondly, such a high traffic rate necessitates the need to use on-chip cache memory on network processors, in order to achieve real-time processing and maintain high throughput. But most on-chip caches on processors are made of SRAM, typically only a few megabytes [3]. A naïve approach that maintains a look-up dictionary to record all the distinct elements seen thus far is too costly, since it requires  $O(n)$  bits. In some applications,  $n$  can be on the order of billions.

Therefore, an ideal algorithm should be simple and quick in processing each flow element and be compact in memory usage. The good news is that it is usually not necessary to know the exact value of a flow’s cardinality. Many approximation algorithms have been developed to explore the trade-off between estimation accuracy and space/timing efficiency.

**Applications Examples:** One important direction of application of cardinality estimation is *network anomaly detection*. Estan et al. [4] suggested three such applications: detecting port scans, detecting denial of service (DoS) attacks and estimating worm spread rate.

Let us take port scan detection as an example to illustrate the idea. A port scan is a probing process that sends service requests from a client to a range of port addresses of a server with the aim of finding an open port, which an attacker can take advantage of to find vulnerabilities on the server. A router can detect such a process by keeping track of packet flows and using cardinality estimation to measure the number of distinct port addresses attempted by each source. Any source that is trying to connect to an abnormally large number of port addresses in a short time interval should be suspected of conducting a port scan.

Other than network anomaly detection, we can apply flow cardinality estimation to many other problems that have a *count-distinct* nature by adapting the concept of a flow to contain other (maybe abstract) types of data. For example, Google might be interested in estimating the number of distinct users that query certain keywords in a period of time [3]. Here a *per-keyword* flow can be defined to be all the user IP addresses that query that keyword. Information about the cardinalities of such flows reflects the popularity of those keywords

---

<sup>2</sup>OC is an acronym for *Optical Carrier*. OC- $n$  specifies the transmission rate of digital signals on a Synchronous Optical Networking (SONET) fiber optic network, equivalent to  $n \times 51.84$  Mbit/s.

which might be helpful in the optimization of Google’s database and search algorithm. Much research on cardinality estimation was motivated by database-related applications (e.g. [5] and [6]).

## 1.2 Motivation and contributions

State-of-the-art algorithms for cardinality estimation comprise two processes. First, the *sketching process* reads elements from the flow and stores useful information in a compact data structure. Second, in the *estimation process*, an estimator takes the recorded information as input and outputs the estimated cardinality of the flow.

The *HyperLogLog* (HLL) algorithm proposed by Flajolet et al. in [7] is near-optimal and has been widely adopted in many industry systems due to its simplicity and excellent performance; according to [3], it is the best existing algorithm for single-flow estimation. Roughly speaking, the HLL algorithm requires hundreds of bytes to make a fair estimation for a single flow with cardinality up to  $4 \times 10^9$ . Details of the HLL algorithm are presented in Section 2.1.1.

However, in some applications the cardinalities of multiple flows need to be estimated at the same time — this is known as per-flow cardinality estimation. In this case, the rate of hundreds of bytes per flow is still too much in some important scenarios where the number of flows can be on the order of tens of millions and memory usage is critical. To this end, Xiao et al. [3] proposed the *virtual LogLog* algorithm based on memory sharing at register level, which can potentially bring down the memory cost from *hundreds of bytes* per flow to *one bit* per flow on average. Details of the sketching process of this algorithm are described in Section 2.2.1. The estimation process of this algorithm uses an estimator called vHLL, summarized in Section 3.2.1.

The main purpose of this thesis is to propose and investigate alternative estimators for the estimation process of the virtual LogLog algorithm in [3]:

- The original vHLL estimator of [3] is based on the HLL estimator for single-flow estimation. We show that the HLL estimator can be generalized by a family of estimators parametrized by a value  $\theta$  (Section 3.1); for HLL,  $\theta = -1$ . The idea of this generalization can be easily applied to the vHLL estimator for per-flow estimation as well. Although it is already known that for single-flow estimation,  $\theta = -1$  (i.e. HLL) is near-optimal [7], it is not clear whether this near-optimality at  $\theta = -1$  extends to the



case of per-flow estimation. We provide empirical evidence to show that indeed  $\theta = -1$  (i.e. vHLL) is still optimal for per-flow estimation.

- We introduce an alternative approach to the estimation problem – a maximum-likelihood estimator. We find that this estimator does not significantly improve the estimation process, compared to the vHLL estimator.

Empirical evidence from both perspectives suggests the near-optimality of the vHLL estimator for per-flow estimation, on the basis of the virtual LogLog algorithm [3].

### 1.3 Thesis overview

The rest of this thesis is organized as follows. Chapter 2 provides the background knowledge and preliminaries to the estimators discussed in this thesis. Chapter 3 and Chapter 4 investigate alternative estimators of the virtual LogLog algorithm from the aforementioned two perspectives. Chapter 3 introduces the  $\text{vLL}_\theta$  estimator, a generalization of the vHLL estimator, and evaluates its performances for different values of  $\theta$ . Chapter 4 derives the maximum-likelihood estimator and compares its performance with that of the  $\text{vLL}_\theta$  estimator. Chapter 5 summarizes the conclusions of the thesis and points to possible future works.

## Chapter 2

# BACKGROUND AND PRELIMINARIES

This chapter briefly overviews cardinality estimation algorithms and points to relevant references. The focus is on summarizing previous works that are essential to the understanding of the main body of this thesis: the LogLog [8] and HyperLogLog [7] algorithms for single-flow estimation and the virtual LogLog algorithm [3] for per-flow estimation.

### 2.1 Single-flow cardinality estimation

Given a single flow of elements  $x_1, x_2, x_3, \dots$ , the goal of cardinality estimation is to estimate the number of distinct elements in the flow. There are two general approaches to this problem: sampling and streaming. A sampling-based algorithm collects a small sample of elements from the flow while bypassing other elements and infers the cardinality of the flow from sample information. A streaming-based algorithm, in contrast, processes every element in the flow. Most state-of-the-art algorithms are based on streaming. For a review of these two kinds of algorithms and a comparison between them, see [9].

A streaming-based algorithm for cardinality estimation has three integral components:

- A *hash function*  $H$ , mapping each element  $x_i$  of the flow into a hashed value  $H(x_i)$  of a chosen type. The purpose of  $H$  is generally twofold: filtering out duplicate elements (relying on the fact that multiple appearances of the same element have the same hash) and providing randomization.
- A *compact data structure*  $\mathcal{D}$ , capturing certain low-dimensional statistical information of the flow's cardinality by recording features of the hashed values. In some literature, the data structure is known as a *sketch*, meaning it is a concise summary of the flow. The process of updating the sketches is called *sketching*.
- An *estimator*  $E$  that takes the content of  $\mathcal{D}$  as input and outputs the estimated cardinality, thus  $\hat{n} = E(\mathcal{D})$ .

Usually we want the estimator to be unbiased (i.e.  $\mathbb{E}[\hat{n}|n] = n$ ). One commonly used metric for the evaluation of the performance of an estimator is the *relative standard error*  $\triangleq \frac{\sqrt{\text{Var}(\hat{n})}}{n}$ . For two unbiased estimators using the same amount of memory, the one with a smaller relative standard error is better. For a detailed classification and comparison of existing single-flow cardinality estimation algorithms, see [10, 11]. We take a closer look at the LogLog and HyperLogLog algorithms in Section 2.1.1.

### 2.1.1 LogLog and HyperLogLog

The HyperLogLog (HLL) [7] algorithm is the best-known algorithm in practice, due to its simplicity and good performance. HLL is a successor of LogLog [8]: the two algorithms use the same hash function and data structure (therefore have the same sketching process) and only differ in their choice of estimators.

The data structure of the LogLog/HLL algorithm is a single array of  $k$  registers (counters), denoted as  $S$ , where  $S[i]$  refers to the  $i^{\text{th}}$  register in the array (or the value stored in it, depending on the context). Suppose we have a hash function  $H$  that maps each element of the flow to a sufficiently long bit string. Define a function  $\rho$  that takes a bit string and outputs the position of the leftmost 1-bit of the string, e.g.  $\rho(000010\cdots) = 5$ . The sketching process of LogLog/HLL is summarized as follows.

#### **LogLog/HLL sketching process:**

1. Initialize all the registers in  $S$  with value 0. Let  $l = \log_2 k$ .
2. For each element  $x$  in flow:
  - (i)  $\langle b_1 b_2 b_3 \cdots \rangle \leftarrow H(x)$ . Hash  $x$  into a sufficiently long bit string.
  - (ii) Divide the bit string into two parts:
    - $j \leftarrow \langle b_1 b_2 \cdots b_l \rangle$ . Use the prefix  $l$  bits to select a register in  $S$ .
    - $q \leftarrow \langle b_{l+1} b_{l+2} \cdots \rangle$ . Use the remaining bits to update the selected register.
  - (iii)  $S[j] \leftarrow \max\{S[j], \rho(q)\}$ . Update the value of  $S[j]$  by the larger of the current value of  $S[j]$  and the position of the leftmost 1-bit of  $q$ .

Here are a few points to note about the above sketching process. First,  $k$  is chosen to be

some integer power of 2 so that  $l$  is an integer and  $j \in \{0, 1, \dots, k-1\}$ . Assume that for a random element  $x$ ,  $H(x)$  is a random bit string with independent and uniform bits. Then  $j$  can be regarded as a uniform random variable in  $\{0, \dots, k-1\}$ , which means that the  $k$  registers in  $S$  are equally likely to be selected for updating for a randomly given element.

Second, the values of the registers in  $S$  after the sketching process should not be affected by duplicate elements, nor by the order of the elements appearing in the flow. This is guaranteed by hashing and the max operation in step (iii).

Third, how large can a register's value get? With the previous assumption on the independence and uniformity of the bits in the hash output string, for a random element  $x$  we have  $\rho(q) = i$  with probability  $\frac{1}{2^i} (i \geq 1)$ . That is,  $\rho(q)$  can be regarded as a geometric random variable with parameter  $1/2$ . By the same assumption, we expect about  $\frac{n}{k}$  distinct elements distributed to any particular register  $S[j]$ . Suppose that this number  $\frac{n}{k}$  is exact, then at the end of the sketching process,  $S[j]$  is just the maximum of  $\frac{n}{k}$  independent  $\text{Geo}(1/2)$  random variables. According to [8], previous study has shown that the expectation of  $S[j]$  is close to  $\log \frac{n}{k}$  with a small additive bias. Therefore, on average each register needs about  $\log \log \frac{n}{k} = O(\log \log n)$  bits to store its value — this is why such a register is also known as a *LogLog sketch*. Based on this discussion, a rough estimator for  $n$  can be  $\hat{n} = k2^{S[j]}$ . Following this idea, the LogLog estimator [8] replaces  $S[j]$  with the average of the  $k$  register values:

$$\hat{n} = \text{LL}(S) \triangleq \eta k 2^{\frac{S[0] + \dots + S[k-1]}{k}}, \quad (2.1)$$

where  $\eta$  is a suitable bias correction factor. The HLL algorithm [7] uses a different estimator:

$$\hat{n} = \text{HLL}(S) \triangleq \gamma k \frac{k}{2^{-S[0]} + \dots + 2^{-S[k-1]}} \quad (2.2)$$

with a different bias correction factor  $\gamma$ . Note that  $2^{\frac{S[0] + \dots + S[k-1]}{k}}$  is commonly known as the *geometric mean* of the values  $2^{S[0]}, \dots, 2^{S[k-1]}$ , while  $\frac{k}{2^{-S[0]} + \dots + 2^{-S[k-1]}}$  is the *harmonic mean* of those values.

We will introduce two other estimators for single-flow estimation in Section 3.1 ( $\text{LL}_\theta$  estimator) and Section 4.4 (maximum-likelihood estimator), respectively.

It has been shown in [8] and [7] that the LogLog and HLL estimators are *asymptotically approximately unbiased* in the sense that, as  $n \rightarrow \infty$ ,  $\frac{\mathbb{E}[\hat{n}]}{n}$  is very close to 1 with a practically negligible fluctuation. It has been shown that as  $n \rightarrow \infty$ , the bias correction constants are independent of  $n$ . They do depend on  $k$ , however; but as  $k$  gets large (e.g.  $k \geq 64$ ), they can be safely replaced by constants for all practical purposes. Practical values of  $\eta$  and  $\gamma$

are 0.39701 and 0.7213 respectively.

The *relative standard error*  $\triangleq \frac{\sqrt{\text{Var}(\hat{n})}}{n}$  is approximately  $\frac{1.30}{\sqrt{k}}$  for the LogLog estimator [8] and approximately  $\frac{1.04}{\sqrt{k}}$  for the HLL estimator [7], as  $n \rightarrow \infty$ . Therefore, using the same number of registers (same amount of memory), HLL’s estimate is more accurate than that of LogLog. With the HLL algorithm, we can achieve an estimate accuracy (in terms of relative standard error) of about 5% by choosing  $k = 512$ . If we use 5 bits for each register (for measuring cardinalities up to  $2^{25} \approx 4 \times 10^9$ ), then in total we need  $5k = 2560$  bits = 320 bytes for one estimation.

According to Section 4 of [7], the HLL algorithm is near-optimal, in the sense that its relative standard error  $\frac{1.04}{\sqrt{k}}$  is quite close to the lower bound  $\frac{1}{\sqrt{k}}$  for a wider class of algorithms based on order statistics. For a more detailed discussion of this lower bound, see [2].

One problem of the HLL algorithm is that it is highly biased and inaccurate for the estimation of small cardinalities. Methods to remedy this flaw are considered in Section 4 of [7] and further in [6].

## 2.2 Per-flow cardinality estimation with memory sharing

In many real-world applications we need to estimate the cardinalities of multiple flows at the same time. For example, consider a stream of data packets from many flows observed by a network monitor device (such as a router). Let each packet be abstracted as a 2-tuple  $(f, x)$ , where  $f$  is the IP source address of the packet (the flow ID) and  $x$  is the destination address (i.e. an element of the flow). The goal is, at the end of the measuring period, to estimate the number of distinct destination addresses from each given source address, i.e. the cardinality of each flow. For example, if the stream of packets is  $(A, 2), (C, 9), (A, 3), (A, 2), (C, 1), (B, 8)$ , then flow  $A$  is  $\{2, 3, 2\}$  with cardinality 2, flow  $B$  is  $\{8\}$  with cardinality 1 and flow  $C$  is  $\{9, 1\}$  with cardinality 2.

An immediate idea to solve this problem is to allocate a separate block of memory for each flow and use any of the existing single-flow cardinality estimation algorithms for each flow. Since we do not know the cardinalities of the flows beforehand, it is inevitable to allocate the maximum amount of memory for each flow; with the HLL algorithm, we still need hundreds of bytes per flow. However, in many applications most of the flows have small cardinalities. Figure 2.1 shows an example of flow cardinality distribution from real-world data.<sup>1</sup> Here, a

---

<sup>1</sup>Data are retrieved from network traffic trace files [12] recorded on a backbone link between San Jose and Los Angeles during a ten minute interval.

flow is defined by an IP source address. The cardinality of a flow is the number of distinct destination addresses among all the packets in the flow. In this example, the majority (about 90%) of the flows have cardinality of only 1, while only six flows have cardinalities larger than  $10^4$ .

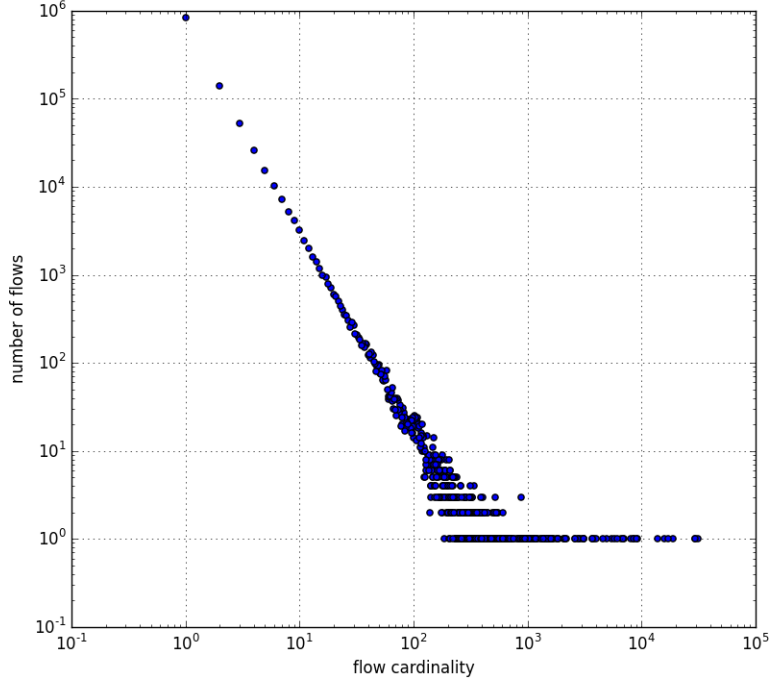


Figure 2.1: A typical flow cardinality distribution from real trace files. x-axis: flow cardinality. y-axis: the number of flows with the corresponding cardinality. Total number of flows: 1,116,535. Average flow cardinality: 2.52.

In view of this waste of memory, some algorithms have been proposed that allow memory to be shared among flows. Such an algorithm usually combines the following components/ideas:

- *A memory pool*, which is the actual source of memory for all estimations.
- *A virtual data structure for each flow*. This data structure does not physically exist but is logically constructed by (often randomly) pointing to memory units in the pool. Consequently, different flows may share some common parts of the memory.
- In the sketching process, for each incoming element, the algorithm applies the sketching process of an existing single-flow estimation algorithm to update the virtual data structure allocated for the corresponding flow.

- After the sketching process, the algorithm uses an estimator to estimate the cardinality of any given flow (with its ID). The estimator needs to take into account the interference among flows due to memory sharing. This estimation process can be done off-line (in comparison to the online sketching process).

Such algorithms may differ from each other in any of the above aspects. A review and comparison of some per-flow cardinality estimation algorithms based on memory sharing can be found in [3]. All these algorithms, however, share memory at *bit* level (meaning that the basic unit of the memory pool and of the virtual data structure is a bit), which [3] claims to be inherently too noisy.

Another algorithm that shares memory at bit level, which [3] did not mention, is the *virtual FM sketches* algorithm proposed in [13]. The algorithm constructs multiple virtual bit arrays from a bit pool for each flow and relies on the sketching process of the FM sketches<sup>2</sup> algorithm for single-flow estimation. The algorithm adopts a maximum-likelihood estimator in the estimation process based on bit patterns in the sketches. We will also consider a maximum-likelihood estimator solution in Chapter 4, but for a different sketching scheme. The performance reported in [13] shows that with memory cost of 1 bit per flow, the algorithm can estimate flows with cardinalities up to 3000 and average flow cardinality about 2.5; it can achieve relative standard error of about 20% for flows with cardinalities less than 500 and 10% for flows with cardinalities in the range 500 to 3000. For relatively small flows with cardinalities in the range 100 to 500, this algorithm may outperform the virtual LogLog algorithm with the vHLL estimator [3]; but it does not appear to be as competitive for larger ranges of flow cardinality.

In view of the inherent drawback of bit-level sharing, [3] proposed the *virtual LogLog* algorithm based on memory sharing at *register* level and showed through experiments that such an algorithm outperforms previous ones. The data structure and sketching process of this algorithm are presented in the following subsection.

### 2.2.1 Virtual LogLog register sharing and sketching process

**Virtual Data Structure:** The virtual LogLog algorithm [3] keeps a memory pool in the form of a register array, denote as  $R$ . Denote the size of the array by  $m$ , which is usually a very large number.  $R[j]$  refers to the  $j^{th}$  register in the array (or the value stored in

---

<sup>2</sup>See [5] by Flajolet and Marin for a description of the FM sketches algorithm for single-flow estimation.

the register, depending on the context). For each flow  $f$ , we form a virtual data structure (*virtual register array*) denoted as  $R_f$ , which is a logically constructed array of  $k$  registers with the  $i^{th}$  register denoted as  $R_f[i]$ . The registers of  $R_f$  are randomly selected from  $R$  by using  $k$  independent hash functions  $G_0, G_1, \dots, G_{k-1}$ , each mapping the flow ID uniformly to an integer in  $\{0, \dots, m-1\}$ , i.e.

$$R_f[i] = R[G_i(f)], \quad 0 \leq i \leq k-1. \quad (2.3)$$

The  $k$  hash functions can be implemented by a single master hash function  $G$  as follows:

$$G_i(f) = G(f|i), \quad (2.4)$$

where “|” is the concatenation operator. It should be emphasized that  $R_f$  does not need to be physically constructed (thus it is “virtual”). A simple example is shown in Figure 2.2 to illustrate the concepts. In the example, say we want to update the third register of flow  $f_2$ , i.e.  $R_{f_2}[3]$ , what actually will be updated is  $R[5]$  — we do not even need to know where the other registers in  $R_{f_2}$  are.

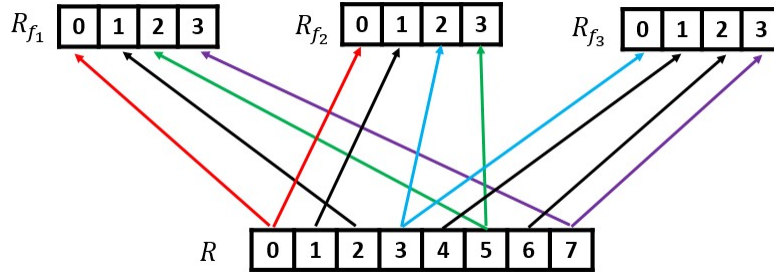


Figure 2.2: An example of the register pool and virtual register arrays. Here  $m = 8$ ,  $k = 4$ . There are three flows and three corresponding virtual register arrays:  $R_{f_1} = [R[0], R[2], R[5], R[7]]$ ,  $R_{f_2} = [R[0], R[1], R[3], R[5]]$ ,  $R_{f_3} = [R[3], R[4], R[6], R[7]]$ .

A caveat should be pointed out here, which the original paper [3] omitted: it could happen that two logically distinct registers of a virtual register array are mapped from the same physical register in the pool, i.e.  $G_i(f) = G_j(f)$  for some  $i \neq j$ . Certainly we wish to avoid this situation, but it is tolerable if the number of such “collisions” is very small.

To see how likely it is for such a “collision” to occur, we can consider a bins-and-balls analogous problem: suppose we have  $k$  balls and  $n$  bins and we throw the balls sequentially, independently and uniformly at random into the bins. In the end how many bins will contain more than one balls? If  $k$  is reasonably large, then the number of balls distributed at *each*



bin can be approximated by an *independent* Poisson random variable  $X$  with mean  $\frac{k}{m}$  (see Chapter 5.4 of [14]). So we have

$$\Pr\{X \geq 2\} = 1 - \Pr\{X \leq 1\} \approx 1 - \left(1 + \frac{k}{m}\right) e^{-\frac{k}{m}} \approx 1 - \left(1 + \frac{k}{m}\right) \left(1 - \frac{k}{m}\right) = \frac{k^2}{m^2},$$

where the last approximation assumes  $k \ll m$ . Therefore out of the  $m$  bins, we expect to have approximately  $\frac{k^2}{m}$  bins that hold more than one ball. Ideally we want  $\frac{k^2}{m}$  to be small, which means with high probability the registers in a virtual register array are all mapped from distinct physical registers in the pool. This factor should be included in the design consideration.

**Sketching Process:** The sketching process of virtual LogLog is *almost* identical to that of LogLog/HLL (recall from Section 2.1.1): for each packet  $(f, x)$  in the stream, we process  $x$  to obtain  $j$  (the register selector) and  $\rho(q)$  (to be compared with the selected register's value); except that the last step (iii) becomes

$$R_f[j] \leftarrow \max\{R_f[j], \rho(q)\}.$$

That is, we are treating the virtual register array as the actual register array. By combining (2.3) and (2.4), the above expression can be re-written as

$$R[G(f|j)] \leftarrow \max\{R[G(f|j)], \rho(q)\}. \quad (2.5)$$

Again, it shows that updates are actually made in the physical registers in  $R$ . It should now be clear why the algorithm is called virtual LogLog: it is based on virtual register arrays, where each register is a LogLog sketch.

After the sketching process, we obtain  $m$  register values. The remaining problem is to infer flow cardinalities from these register values: this is the *estimation process*. We will consider two kinds of estimators in Chapters 3 and 4, respectively.

## 2.3 A per-flow estimator performance metric – Weighted square error

Before we discuss specific estimators for per-flow cardinality estimation, we consider a metric based on weighted square error to evaluate the performance of any given estimator.

Suppose that the incoming data stream contains  $M$  flows, with cardinalities  $n_1, n_2, \dots, n_M$ . Given any estimator  $E$ , suppose its corresponding estimates for the flows' cardinalities are  $\hat{n}_1, \hat{n}_2, \dots, \hat{n}_M$ . The *weighted square error* of estimator  $E$  is defined as

$$\text{WSE}(E) \triangleq \sum_{i=1}^M (n_i - \hat{n}_i)^2 w(n_i), \quad (2.6)$$

where  $w(\cdot)$  is a weight function mapping the cardinality of a flow to a positive number. For two estimators using the same amount of memory, the one with a lower WSE is better.

### 2.3.1 Weight function

The choice of the weight function depends on the specific application. For example, if each flow is considered equally important regardless of its cardinality, then we can simply let  $w(n) = 1$  for all  $n$ . In many other situations, large flows are considered to be more important than small flows, then we want  $w(n)$  be an increasing function in  $n$ . The weight function used in this thesis is presented and explained as follows.

First, assume that the cardinality of a randomly chosen flow can be modeled as a random variable  $N$  with pmf  $p_N(n) \triangleq \Pr(N = n)$ . We use the following weight function:

$$w(n) = \frac{1}{n \cdot p_N(n)}. \quad (2.7)$$

The motivation for using this weight function is explained here. The integral

$$\int_{m_1}^{m_2} w(n) p_N(n) \, dn$$

is an approximation of the total weight of flows whose cardinalities are in the range  $[m_1, m_2)$ . With the weight function in (2.7), we have  $w(n)p_N(n) = \frac{1}{n}$ ; it is easy to verify that in this case the integral

$$\int_{n^*}^{n^*(1+\epsilon)} w(n) p_N(n) \, dn$$

for  $\epsilon > 0$  is independent of  $n^*$ . If we plot the curve  $w(n)p_N(n)$  as a function of  $n$  on log scale of  $n$ , then the area under the curve should be approximately the same in each decade interval:  $[1, 10)$ ,  $[10, 100)$ ,  $[100, 1000)$ ,  $[1000, 10000)$ , etc. In another word, by choosing such a weight function, we put approximately the same total weight to the aggregate of flow

cardinalities in each of these intervals.

### 2.3.2 Zipf model

The weight function in (2.7) can be applied to any given flow cardinality distribution. The particular distribution used for simulation in this thesis is the Zipf distribution. A random variable  $N$  following the  $\text{Zipf}(\pi, n_{\max})$  distribution has the following pmf:

$$p_N(n) = \frac{n^{-\pi}}{C}, \quad n \in \{1, 2, \dots, n_{\max}\}, \quad (2.8)$$

where  $n_{\max}$  is an upper bound of the flow cardinality,  $\pi$  is a parameter that controls the shape of the Zipf distribution ( $\pi > 0$ ) and  $C$  is a constant such that  $\sum_{n=1}^{n_{\max}} p_N(n) = 1$ .

The adoption of the Zipf model is motivated by the fact that Zipf's law underlies many Internet applications (see [15]). In Figure 2.3, we plot the distribution of a large number of simulated flow cardinalities, each generated independently according to a  $\text{Zipf}(2.25, 10^5)$  model. The validity of the model can be verified by observing the resemblance between Figure 2.3 (from Zipf model) and Figure 2.1 (from raw Internet data).

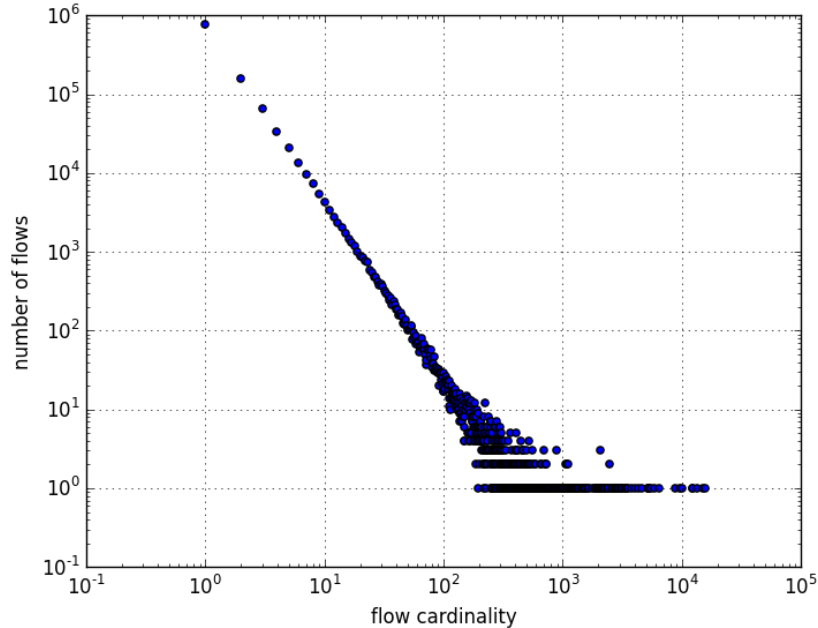


Figure 2.3: Distribution of simulated flow cardinalities according to  $\text{Zipf}(2.25, 10^5)$ . Total number of flows: 1,116,535. Average flow cardinality: 2.88.

With  $N \sim \text{Zipf}(\pi, n_{\max})$ , by (2.7) and (2.8), we have

$$w(n) = Cn^{\pi-1}, \quad (2.9)$$

where the constant  $C$  can be omitted (set to 1) because when we compare the weighted square errors of two estimators, we only care about their relative values, which are unaffected by a constant factor.

We remark that the specific estimators we will discuss in Chapters 3 and 4 do not rely on any particular distribution model. The Zipf model here is for two purposes: to complete the definition of the weighted square error for the performance evaluation of estimators and to generate random flow cardinalities for simulation; both purposes are independent of the estimation process.

Finally, we make a note on experimental evaluation of per-flow estimators. Given the close match between the real trace data and simulated trace data with a Zipf distribution, we will evaluate the estimators using simulated trace data. We generated 100 simulated trace files. Each trace file contains  $10^6$  flows, with the cardinality of each flow randomly and independently generated according to a  $\text{Zipf}(2.25, 10^5)$  distribution. A flow is presented in the file as a collection of packets with the same source address and distinct destination addresses; different flows have different source addresses (flow IDs). For one experiment, we process one such trace file and estimate the cardinalities of all the  $10^6$  flows in it. For a given estimator, we perform 100 independent experiments and evaluate its performance based on the statistics averaged over the 100 experiments. All experiments are performed with  $m = 200,000$  and  $k = 512$ . If each register uses 5 bits, this setting uses  $10^6$  bits to measure the cardinalities of  $10^6$  flows, leading to one bit per flow on average.

## Chapter 3

# VIRTUAL LOGLOG ESTIMATOR WITH PARAMETER $\theta$

The sketching process of the virtual LogLog algorithm has been described in Section 2.2.1. After the sketching process, we can offload the register values from the network measuring device for off-line query. The following per-flow estimation problem is to be solved:

**Given:** Register values  $R[0], R[1], \dots, R[m-1]$  and any flow's ID  $f$ .

**Objective:** Estimate  $n_f$ , the number of distinct elements in flow  $f$ .

In this chapter, we first introduce  $LL_\theta$ , a family of generalized LogLog estimators parameterized by  $\theta$ , for single-flow estimation. Then, with the similar idea from  $LL_\theta$ , we propose  $vLL_\theta$ , a family of generalized virtual LogLog estimators parameterized by  $\theta$ , for per-flow estimation.

### 3.1 The $LL_\theta$ estimator for single-flow estimation

The  $LL_\theta$  estimator is a class of estimators parameterized by  $\theta$ , which unifies and generalizes the LogLog and HLL estimators described in Section 2.1.1.

#### 3.1.1 Generalized mean

We start with the concept of *generalized mean*. Given  $k$  positive numbers  $x_1, x_2, \dots, x_k$ , their *generalized mean parameterized by  $\theta$*  is defined for nonzero  $\theta$  by

$$A_\theta(x_1, \dots, x_k) = \left( \frac{1}{k} \sum_{i=1}^k x_i^\theta \right)^{\frac{1}{\theta}}, \quad \theta \in \mathbb{R}, \theta \neq 0. \quad (3.1)$$

In the case of  $\theta = 0$ , we let

$$A_0(x_1, \dots, x_k) = \left( \prod_{i=1}^k x_i \right)^{\frac{1}{k}}, \quad (3.2)$$

which is in fact the limit of  $A_\theta(x_1, \dots, x_k)$  as  $\theta \rightarrow 0$ . Note that

$$A_\theta(x_1, \dots, x_k) = \begin{cases} \min\{x_1, \dots, x_k\} & \text{if } \theta = -\infty \\ \frac{k}{\frac{1}{x_1} + \dots + \frac{1}{x_k}} & \text{if } \theta = -1 \\ \frac{x_1 + \dots + x_k}{k} & \text{if } \theta = 1 \\ \max\{x_1, \dots, x_k\} & \text{if } \theta = \infty. \end{cases} \quad (3.3)$$

$A_\theta(x_1, \dots, x_k)$  is commonly known as the *arithmetic mean*, *geometric mean*, and *harmonic mean* of the  $k$  numbers when  $\theta = 1, 0, -1$ , respectively. A notable property of the generalized mean function is that  $A_\theta(x_1, \dots, x_k)$  with a lesser  $\theta$  is more robust to abnormally high values in obtaining the mean. Consider an example: for numbers 1, 1, 1, 1, 100,  $A_1 = 20.8$ ,  $A_0 \approx 2.51$ ,  $A_{-1} \approx 1.25$  and  $A_{-\infty} = 1$ .

### 3.1.2 Unification and generalization of LogLog and HLL estimators

Recall from Section 2.1.1, the LogLog and HLL estimators use the geometric mean and harmonic mean, respectively. Based on the generalized mean notation, we attempt to unify these two estimators by proposing the  $LL_\theta$  estimator:

$$\hat{n} = LL_\theta(S) \triangleq \xi k A_\theta(2^{S^{[0]}}, \dots, 2^{S^{[k-1]}}), \quad (3.4)$$

where  $\xi$  is a suitable coefficient.

It is desirable if there exists a value of  $\xi$  to make the estimator approximately unbiased and for which we can identify such value. Specifically,  $\xi$  should not depend on  $n$ . We already know that for  $\theta = 0$  (i.e. LogLog [8]) and  $\theta = -1$  (i.e. HLL [7]) such values of  $\xi$  exist (by letting  $\xi = \eta$  and  $\xi = \gamma$  respectively). But we do not know if it is true for other values of  $\theta$ . The analysis of this estimator for general  $\theta$  is difficult.<sup>1</sup> We resort to simulations to explore empirical evidence of the existence of this coefficient  $\xi$ .

---

<sup>1</sup>Interested readers may refer to the analyses of LogLog [8] and HLL [7] algorithms to get an idea of the techniques used for this kind of analysis.

Let  $A_\theta$  denote  $A_\theta(2^{S[0]}, \dots, 2^{S[k-1]})$  for short. In Figure 3.1 we show empirical values of the ratio  $\frac{n}{kA_\theta}$  for selected values of  $k$ ,  $n$  and  $\theta$ . Each dot in each of the sub-figures is generated by taking the average of 50 independent experiment results. Experiments are performed on values of  $\theta$  in  $[-3.0, -2.9, -2.8, \dots, 0.9, 1.0]$ , but plots are only shown for selected values of  $\theta$  for better graph layout; plots for other values of  $\theta$  have similar shapes and are at their expected positions in the figure.

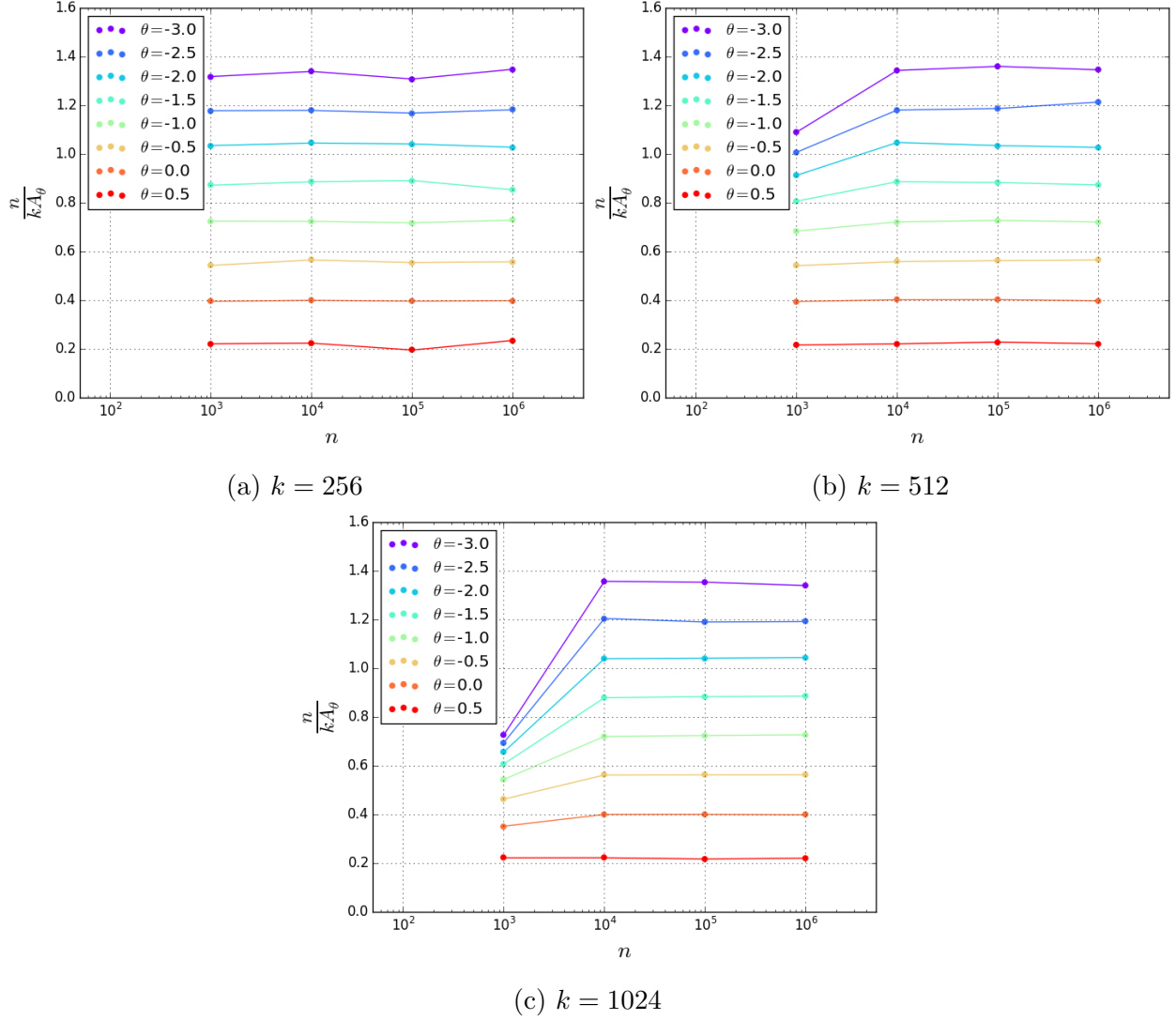


Figure 3.1: Empirical values of  $\frac{n}{kA_\theta}$  for different  $k$ ,  $n$  and  $\theta$ 's.

From the plots we see:

- For the same value of  $\theta$  and  $k$ ,  $\frac{n}{kA_\theta}$  is almost a constant when  $n$  gets large (say  $n > 10000$ ).

- For the same value of  $\theta$  and  $n$ ,  $\frac{n}{kA_\theta}$  is almost a constant for large  $k$  (256, 512 or 1024).

We conclude that, at least for  $\theta \in [-3.0, 1.0]$ , there exists a value for the aforementioned coefficient  $\xi$  that approximately depends only on  $\theta$  for large  $n$  and  $k$ . To reflect this dependence of  $\xi$  on  $\theta$ , we denote this coefficient as  $\xi_\theta$  instead.

To investigate the relationship between  $\xi_\theta$  and  $\theta$ , we plot values of  $\xi_\theta$  computed empirically for selected values of  $\theta$  in  $[-3.0, 1.0]$ , shown in Figure 3.2. Here, the values are generated empirically with fixed  $n = 10^5$  and  $k = 512$ .

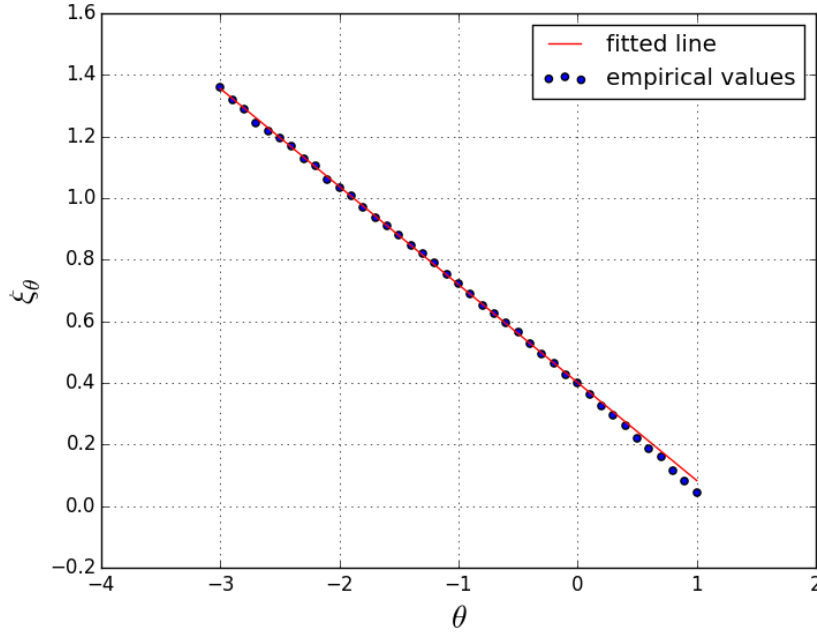


Figure 3.2: Empirical values of  $\xi_\theta$  as a function of  $\theta$  and a fitted line showing the approximate linear relationship.

We observe an approximately linear relationship between  $\xi_\theta$  and  $\theta$  by

$$\xi_\theta = 0.401 - 0.318\theta, \quad -3.0 \leq \theta \leq 1.0, \quad (3.5)$$

which is plotted as the red line in Figure 3.2. The  $LL_\theta$  estimator is now completely defined by replacing  $\xi$  in (3.4) with  $\xi_\theta$  specified in (3.5). The claim that this family of estimators unifies the LogLog and HLL estimators can be verified by checking that  $\xi_0 \approx \eta$  (for LogLog) and  $\xi_{-1} \approx \gamma$  (for HLL).

A natural question one may ask next is if there exists an optimal value of  $\theta$  for the  $LL_\theta$  estimator. To answer this question, we empirically calculate the relative standard error of



the  $LL_\theta$  estimator for different values of  $\theta$  (and for selected  $k$ 's). The results are plotted in Figure 3.3.

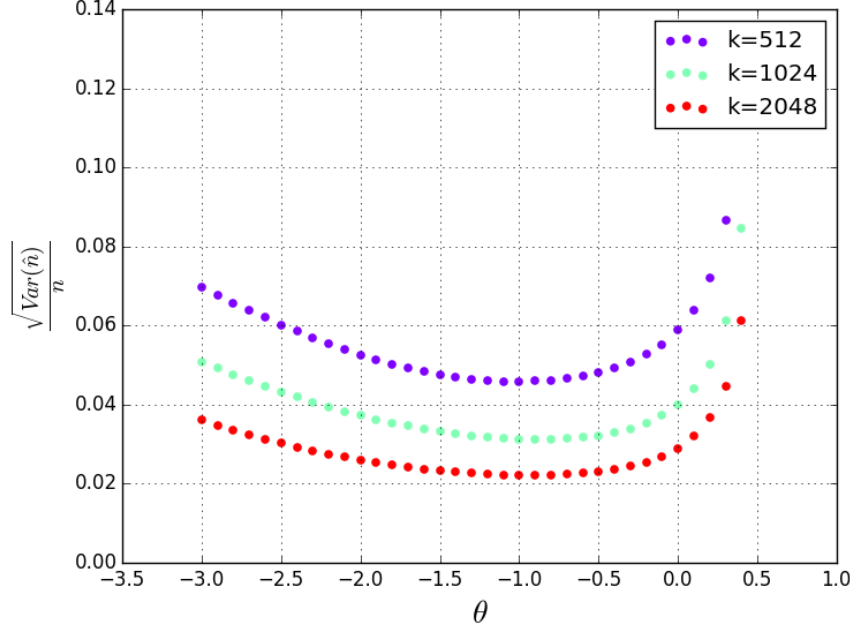


Figure 3.3: Empirically calculated relative standard error  $\triangleq \frac{\sqrt{\text{Var}(\hat{n})}}{n}$  of the  $LL_\theta$  estimator vs.  $\theta$  for selected  $k$ 's.

According to the experiment results, for each fixed value of  $k$ , the optimal value of  $\theta$  for the  $LL_\theta$  estimator is either  $-1$  or  $-0.9$ ; in the cases where  $LL_{-0.9}$  is the optimal, the difference between  $LL_{-1}$  and  $LL_{-0.9}$ , in terms of relative standard error, is negligible. Recall the  $LL_{-1}$  estimator is the same as the HLL estimator. This verifies the claim that the HLL is near-optimal [7], which we briefly described in Section 2.1.1.

### 3.2 The $vLL_\theta$ estimator for per-flow estimation

With the  $LL_\theta$  estimator defined, we now introduce the  $vLL_\theta$  estimator for per-flow estimation. We start with a high-level idea of the estimator and the motivation for introducing the parameter  $\theta$  here.

Just like for single-flow estimation where we infer the cardinality of a flow from its register array, for per-flow estimation we can infer the cardinality of a given flow  $f$  from its *virtual*

register array  $R_f$ . We can directly apply the  $LL_\theta$  estimator on  $R_f$  to give an estimate of the total number of distinct elements distributed to  $R_f$ . But this estimate involves the noise brought by other flows that have registers shared with flow  $f$ ; so we probably need a different bias correction coefficient to calibrate this rough estimate. While for the single-flow case  $\theta = -1$  is near-optimal, it is not immediately clear whether  $\theta = -1$  with a suitable bias correction coefficient is still near-optimal in the per-flow case.

Recall that  $A_\theta$  with a lesser value of  $\theta$  is more robust to abnormally high values in obtaining the mean. For per-flow estimation, due to register sharing, large flows can cause much noise to some small flows by causing abnormally high register values. Therefore, we speculate that an estimator with a lesser value of  $\theta$  might work better in the per-flow case due to its robustness against bursty noise from large flows. This motivates our investigation of the  $vLL_\theta$  estimator, introduced in the rest of this section.

### 3.2.1 Review of vHLL estimator and the generalization to $vLL_\theta$ estimator

We start with a review of the virtual HyperLogLog (vHLL) estimator introduced in [3]. Let  $n_T$  be the total aggregate cardinality of all flows, i.e. the sum of the cardinalities of all flows in the packet stream. Let  $n_{f+}$  be the total number of distinct elements distributed to the register array  $R_f$ , which include the distinct elements from flow  $f$  and those from other flows (we call noise elements). Suggested in [3] is the following approximate relationship between  $n_f, n_{f+}$  and  $n_T$ :

$$n_{f+} - n_f \approx \frac{k}{m}(n_T - n_f), \quad (3.6)$$

which can be interpreted as: the number of noise elements received by  $R_f$  is the total number of elements from flows other than  $f$  scaled by  $\frac{k}{m}$  — the ratio of the number of registers in  $R_f$  to the number of registers in  $R$ . The assumption here is that noise elements are roughly uniformly distributed to all the  $m$  registers in the pool, which is a fair approximation when the number of flows and the number of registers for each flow are both sufficiently large [3]. Rearranging (3.6) gives us

$$n_f \approx -\frac{k}{m-k}n_T + \frac{m}{m-k}n_{f+}. \quad (3.7)$$

Now the values of  $n_T$  and  $n_{f+}$  are not directly available so we need their estimates  $\hat{n}_T$  and  $\hat{n}_{f+}$ . It has been shown in Section 6.1 of [3] that if  $\hat{n}_{f+}$  and  $\hat{n}_T$  are close to  $n_f + \frac{k}{m}(n_T - n_f)$  (which approximately equals to  $n_{f+}$  by (3.6)) and  $n_T$  respectively, then  $\hat{n}_f$  is an approximately

unbiased estimator of  $n_f$ .

There are two possible methods to estimate  $n_T$ . The first method is to treat  $n_T$  as the cardinality of a *grand flow* – the flow containing all the distinct packets in the stream. In another word, if a packet is abstracted as a  $(f, x)$  pair, to estimate  $n_T$  is to estimate the total number of distinct  $f - x$  pairs (e.g. source-destination pair) in the stream. This is a single-flow cardinality estimation problem and, as discussed in Section 2.1, can be solved using the HLL algorithm with an additional few hundreds of bytes, which is negligible compared to the memory for main register pool  $R$ . The second method is to use  $\hat{n}_T = \text{HLL}(R)$  as a rough estimator, based on the assumption that all the elements of the grand flow are distributed approximately uniformly over  $R$ . Either method works fine in practice. So  $\hat{n}_T$  is relatively easy to obtain.

Now, for  $n_{f+}$ , the vHLL estimator applies the HLL single-flow estimator on  $R_f$  to obtain an estimate, i.e.

$$\hat{n}_{f+} = \text{HLL}(R_f) = \xi_{-1} k A_{-1} (2^{-R_f[0]}, \dots, 2^{-R_f[k-1]}). \quad (3.8)$$

Combining the ideas above, the vHLL estimator for flow  $f$  is summarized as:

$$\hat{n}_f = -\frac{k}{m-k} \hat{n}_T + \frac{m}{m-k} \xi_{-1} k A_{-1} (2^{-R_f[0]}, \dots, 2^{-R_f[k-1]}). \quad (3.9)$$

We consider a generalization of the vHLL estimator by replacing the above harmonic mean  $A_{-1}$  with a generalized mean  $A_\theta$ . More specifically, given any value  $\theta$ , let

$$\tilde{n}_f(\theta) = k \cdot A_\theta (2^{R_f[0]}, \dots, 2^{R_f[k-1]}) \quad (3.10)$$

be a rough estimate of  $n_f$ . Then we calibrate this rough estimate by suitable additive and multiplicative constants  $\alpha_\theta$  and  $\beta_\theta$  to obtain a better estimate, i.e.

$$\hat{n}_f(\theta) = \alpha_\theta + \beta_\theta \tilde{n}_f(\theta). \quad (3.11)$$

We discuss how to set the values of  $\alpha_\theta$  and  $\beta_\theta$  in Section 3.2.2.

### 3.2.2 How to set $\alpha_\theta$ and $\beta_\theta$

For a given  $\theta$ , one way to set the values of  $\alpha_\theta$  and  $\beta_\theta$  is through *empirical error minimization*, explained as follows. Recall from the end of Section 2.3.2, we have 100 simulated trace files, each can be considered as a training sample. Suppose that in one training sample we have  $M$  flows with cardinalities  $n_1, n_2, \dots, n_M$ , and by the vLL $_\theta$  estimator specified in (3.11) we obtain corresponding estimates  $\hat{n}_1(\theta), \hat{n}_2(\theta), \dots, \hat{n}_M(\theta)$ . Then we can find values of  $\alpha_\theta$  and  $\beta_\theta$  that minimize the *weighted square error* defined in Section 2.3 on this training sample. That is

$$(\alpha_\theta, \beta_\theta) = \arg \min_{(\alpha, \beta)} \sum_{i=1}^M (n_i - \hat{n}_i(\theta))^2 w(n_i) = \arg \min_{(\alpha, \beta)} \sum_{i=1}^M (n_i - (\alpha + \beta \tilde{n}_i(\theta)))^2 w(n_i), \quad (3.12)$$

where we let  $w(n_i) = n_i^{\pi-1}$  and  $\pi = 2.25$ . The solution to the above minimization problem is standard:

$$\beta_\theta = \frac{\overline{xy} - \bar{x} \cdot \bar{y}}{\overline{x^2} - \bar{x}^2}, \quad \alpha_\theta = \bar{y} - \beta_\theta \bar{x}, \quad (3.13)$$

where

$$\bar{x} = \frac{\sum_{i=1}^M w(n_i) \tilde{n}_i(\theta)}{\sum_{i=1}^M w(n_i)}, \bar{y} = \frac{\sum_{i=1}^M w(n_i) n_i}{\sum_{i=1}^M w(n_i)}, \overline{xy} = \frac{\sum_{i=1}^M w(n_i) \tilde{n}_i(\theta) n_i}{\sum_{i=1}^M w(n_i)}, \overline{x^2} = \frac{\sum_{i=1}^M w(n_i) \tilde{n}_i^2(\theta)}{\sum_{i=1}^M w(n_i)}.$$

Since we have 100 training samples, we can obtain the values of  $\alpha_\theta$  and  $\beta_\theta$  for each of the samples using the above method and use their average values for the estimator. The values of  $\alpha_\theta$  and  $\beta_\theta$  obtained in this way are optimal in the sense that they minimize the weighted square error of the training samples (i.e. the *empirical error*). The drawback of this approach is obvious. It only works well for the training samples generated according to a specific Zipf distribution; there is no performance guarantee for other data sets. Also, it relies on the specific definition of the weighted square error (including the weight function and the cardinality distribution model), which may not be universal for all applications.

We suggest the following alternative and practical choice of  $\alpha_\theta$  and  $\beta_\theta$  for the vLL $_\theta$  estimator, which does not rely on any assumption on weight function or cardinality distribution model:

$$\beta_\theta = \frac{m}{m-k} \xi_\theta \quad \text{and} \quad \alpha_\theta = -\frac{k}{m-k} \hat{n}_T, \quad (3.14)$$

where  $\xi_\theta$  can be calculated using (3.5) and  $\hat{n}_T$  can be calculated using either of the two methods mentioned in Section 3.2.1 (e.g.  $\hat{n}_T = \text{HLL}(R)$ ). Note that this is a direct generalization

of the vHLL estimator in (3.9). Also note that in this case  $\alpha_\theta$  does not depend on  $\theta$ .

In Figure 3.4, we plot the empirical values of  $\alpha_\theta$  and  $\beta_\theta$  obtained by these two different methods through experiments on the 100 training samples. The plots show that our suggested values of  $\alpha_\theta$  and  $\beta_\theta$  are good as they are close to the optimal values of  $\alpha_\theta$  and  $\beta_\theta$  for these training samples.

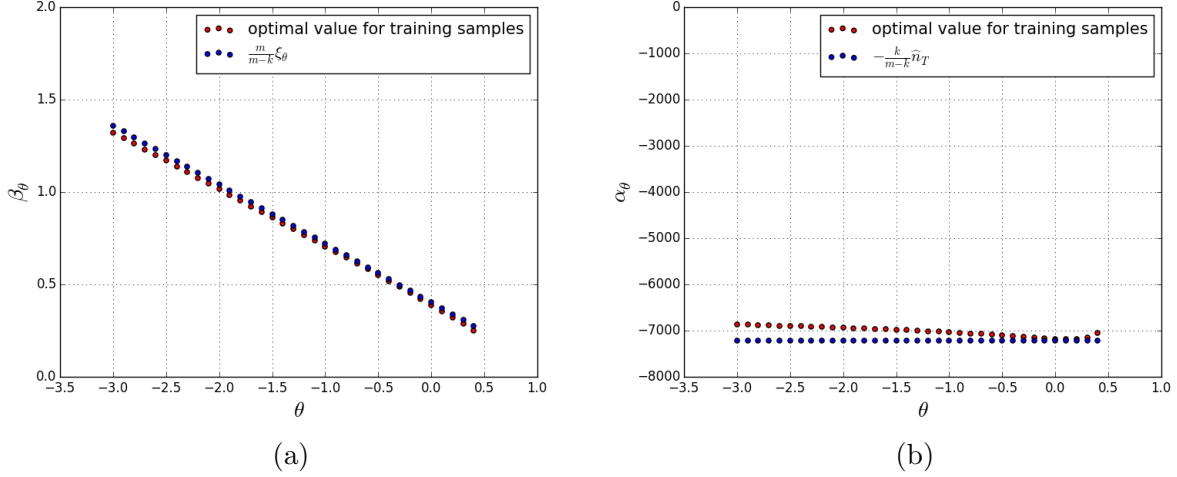


Figure 3.4: The values of  $\alpha_\theta$  and  $\beta_\theta$  vs.  $\theta$  by two different methods.

### 3.2.3 Summary of the vLL $_\theta$ estimator

We summarize the vLL $_\theta$  estimator here. For a given  $\theta \in [-3.0, 1.0]$ , the vLL $_\theta$  estimator estimates the cardinality of a given flow  $f$  by

$$\hat{n}_f = -\frac{k}{m-k}\hat{n}_T + \frac{m}{m-k}\xi_\theta k A_\theta(2^{R_f[0]}, \dots, 2^{R_f[k-1]}), \quad (3.15)$$

where

$$\xi_\theta = 0.401 - 0.318\theta, \quad -3.0 \leq \theta \leq 1.0, \quad (3.16)$$

and  $\hat{n}_T$  is an estimate of the total aggregate cardinality of all the flows and can be obtained by calculating HLL( $R$ ) in practice.

Note that it could happen that Equation (3.15) produces an estimate  $\hat{n}_f < 1$ , which is impossible in practice. In such case, we can simply reset  $\hat{n}_f = 1$ .

### 3.3 Experimental performance evaluation

We ran experiments on the 100 simulated trace files to evaluate the performance of the  $\text{vLL}_\theta$  estimator for different values of  $\theta$ , in terms of the weighted square error defined in Section 2.3. The results are plotted in Figure 3.5. Results for  $\theta \geq 0.5$  are not plotted in the graph because the values are comparatively too large (which means the corresponding  $\text{vLL}_\theta$  estimators are bad).

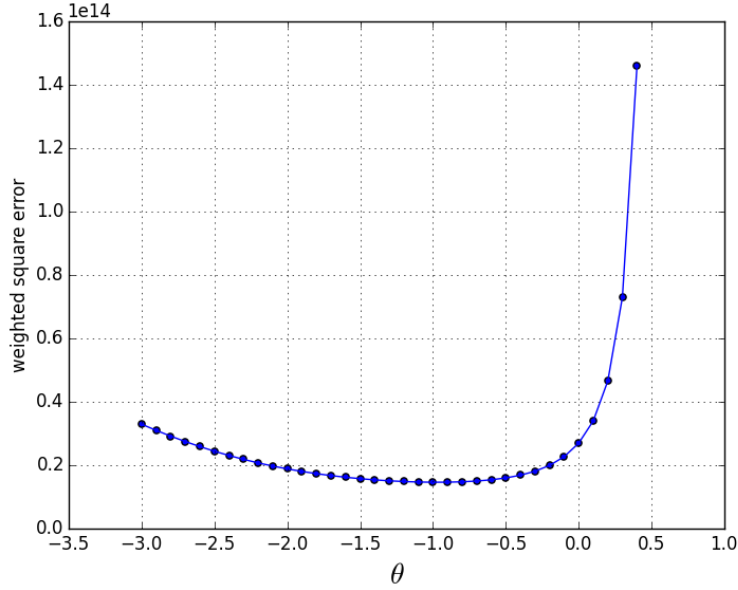


Figure 3.5: Experimental values of  $\text{WSE}(\text{vLL}_\theta)$  vs.  $\theta$ .

According to the experiment results, the  $\text{vLL}_{-1}$  estimator has the best performance, which refutes our speculation at the beginning of Section 3.2 that a value of  $\theta$  lesser than  $-1$  might be optimal for the per-flow estimation case.

A possible explanation of this phenomenon is as follows. With the weight function defined in Section 2.3 and  $\text{Zipf}(2.25, 10^5)$  used for flow cardinality distribution, we have that  $w(n) \propto n^{1.15}$ . This means we put much larger weights on large flows compared to small flows. Therefore the weighted square error of an estimator is largely determined by how well the estimator estimates large flows. But the influence of bursty noise on large flows is much weaker than that on small flows; in another word, the signal-to-noise ratio for the cardinality estimation of a large flow is larger than that for small flows. Therefore, as we have discussed for the single-flow case where  $\theta = -1$  is near-optimal, for per-flow case  $\theta = -1$  is also likely to give good estimates for large flows and hence results in a smaller weighted square error.

## Chapter 4

# MAXIMUM LIKELIHOOD ESTIMATOR

In this chapter we derive an alternative approach to the same per-flow estimation problem stated at the beginning of Chapter 3: a maximum-likelihood estimator.

### 4.1 Formulation

We model the values of the  $k$  registers in  $R_f$  after the sketching process as a random vector

$$Z_f = (Z_{f,0}, Z_{f,1}, \dots, Z_{f,k-1}),$$

where  $Z_{f,i} \in \{0, \dots, r_{\max}\}$ . The  $r_{\max}$  is the maximum value that can be stored in a register. For example, if the register has 5 bits, then  $r_{\max} = 2^5 - 1 = 31$ . Register values  $R_f = (R_f[0], \dots, R_f[k-1])$  is an instance of the random vector  $Z_f$ . Let  $p_n(z_0, \dots, z_{k-1})$  be the pmf of  $Z_f$  (i.e. the joint pmf of the  $k$  random variables) given that  $n_f = n$ , i.e. it is the likelihood function of  $n$ :

$$p_n(z_0, \dots, z_{k-1}) \triangleq \Pr \left\{ Z_f = (z_0, \dots, z_{k-1}) \mid n_f = n \right\}. \quad (4.1)$$

The maximum likelihood estimator finds the value of  $n$  that maximizes the likelihood function based on the instance of register values in  $R_f$ . An equivalent formulation is to maximize the natural log of the likelihood function

$$L_f(n) \triangleq \ln p_n(R_f[0], \dots, R_f[k-1]). \quad (4.2)$$

The maximum-likelihood estimator of flow  $f$ 's cardinality is then given by

$$\hat{n}_{f,\text{ML}} = \arg \max_n L_f(n), \quad n \in \{1, 2, \dots, n_{\max}\}, \quad (4.3)$$

where  $n_{\max}$  is an upper bound of the cardinality of a flow.

## 4.2 Derivation of the likelihood function

In this section we derive an analytical expression for the likelihood function  $p_n(z_0, \dots, z_{k-1})$ .

### 4.2.1 Outline of derivation

Suppose we have two copies of the register pool  $R$ :  $R^{(1)}$  and  $R^{(2)}$ , both initialized with 0's. Consider the following process:

1. Repeat the virtual LogLog sketching process described in Section 2.2.1 for all elements, except those from flow  $f$ , with register pool  $R^{(1)}$ .
2. Repeat the sketching process only for elements from flow  $f$  with register pool  $R^{(2)}$ .
3. Merge the register values in  $R^{(1)}$  and  $R^{(2)}$  by a register-wise max operation.

The register values in the merged pool as described above are exactly the same as those in  $R$  after the original sketching process. This is because the register values at the end of the sketching process are not affected by the order of arrival of the elements in the stream.

Let  $R_f^{(1)}$  ( $R_f^{(2)}$ ) denote the virtual register array for flow  $f$  constructed by selecting registers from  $R^{(1)}$  ( $R^{(2)}$ ), whose indices in  $R^{(1)}$  ( $R^{(2)}$ ) are the same as those of  $R_f$  in  $R$ . Then define the following random variables:

$Z' \triangleq$  the value of an arbitrary register in  $R^{(1)}$ .

$W_n \triangleq$  the value of an arbitrary register in  $R_f^{(2)}$ , conditioned on  $n_f = n$ .

$Z_n \triangleq$  the value of an arbitrary register in  $R_f$ , conditioned on  $n_f = n$ .

$Z'$  represents the background noise in our estimation caused by elements from flows other than  $f$ .  $W_n$  represents the impact of elements from flow  $f$  itself. We model the values of the  $k$  corresponding registers in  $R_f^{(1)}$  and  $R_f^{(2)}$  also as random vectors, respectively, i.e.

$$\begin{aligned} Z_f^{(1)} &= (Z_{f,0}^{(1)}, Z_{f,1}^{(1)}, \dots, Z_{f,k-1}^{(1)}), \\ Z_f^{(2)} &= (Z_{f,0}^{(2)}, Z_{f,1}^{(2)}, \dots, Z_{f,k-1}^{(2)}). \end{aligned}$$



Then we have

$$\begin{aligned}
Z_{f,0}^{(1)}, Z_{f,1}^{(1)}, \dots, Z_{f,k-1}^{(1)} &\sim Z', \\
Z_{f,0}^{(2)}, Z_{f,1}^{(2)}, \dots, Z_{f,k-1}^{(2)} &\sim W_n, \\
Z_{f,0}, Z_{f,1}, \dots, Z_{f,k-1} &\sim Z_n, \\
Z_{f,i} &= \max\{Z_{f,i}^{(1)}, Z_{f,i}^{(2)}\}, \quad i = 0, 1, \dots, k-1.
\end{aligned}$$

Therefore

$$Z_n = \max\{Z', W_n\}. \quad (4.4)$$

Here is an outline of the derivation of an analytical expression for  $p_n(z_0, \dots, z_{k-1})$ :

1. Find the distribution of  $Z'$  and argue that  $Z_{f,0}^{(1)}, Z_{f,1}^{(1)}, \dots, Z_{f,k-1}^{(1)}$  are independent. We will see that it is difficult to obtain the exact distribution of  $Z'$ , but there exists a convenient and close approximation to it.
2. Find the distribution of  $W_n$  and argue that  $Z_{f,0}^{(2)}, Z_{f,1}^{(2)}, \dots, Z_{f,k-1}^{(2)}$  are independent under the assumption that  $n$  is large.
3. With the distribution of  $Z'$  and  $W_n$ , find the distribution of  $Z_n$  by (4.4).
4. With the independence of  $Z_{f,0}^{(1)}, Z_{f,1}^{(1)}, \dots, Z_{f,k-1}^{(1)}$  and  $Z_{f,0}^{(2)}, Z_{f,1}^{(2)}, \dots, Z_{f,k-1}^{(2)}$ , we have that  $Z_{f,0}, Z_{f,1}, \dots, Z_{f,k-1}$  are independent. Then we can factor  $p_n(z_0, \dots, z_{k-1})$  out as the product of  $k$  pmfs of a single variable  $Z_n$ .

#### 4.2.2 Distribution of $Z'$ and independence of $Z_{f,0}^{(1)}, \dots, Z_{f,k-1}^{(1)}$

One possible approach to obtain an approximate expression of the distribution of  $Z'$  is as follows. First assume that the total number of flows (excluding flow  $f$ ) is known and a flow's cardinality can be modeled as a random variable (e.g. Zipf). Then we can model the number of distinct elements distributed to an arbitrary register by a random variable  $Y$  and calculate the distribution of  $Y$  by adding up the influences from all the individual independent flows on that register. This calculation involves a high-dimensional convolution, which is computationally costly. We can perhaps avoid the convolution by appealing to the central limit theory for approximation. However this approximation gives us a continuous distribution for  $Y$ , from which we need to derive the CDF of  $Z'$  that only takes discrete values (and really concentrates on only a few values as we shall see soon). This calculation

may be highly inaccurate and is somewhat complicated. Moreover, in practice we usually do not know the total number of flows beforehand and it is undesirable that our estimator relies on any particular cardinality distribution model.

Another approach to approximate the distribution of  $Z'$ , which circumvents the above difficulties and complications, is the following. Define a random variable  $Z$ :

$$Z \triangleq \text{the value of an arbitrary register in } R.$$

We claim that, for any practical purpose, the distribution of  $Z'$  can be directly approximated by the distribution of  $Z$ , i.e.

$$F_{Z'}(i) \approx F_Z(i), \quad 0 \leq i \leq r_{\max}, \quad (4.5)$$

where  $F_{Z'}$  and  $F_Z$  are the CDFs of  $Z'$  and  $Z$ , respectively. This approximation is based on the following two considerations:

- Register values in  $R^{(1)}$  differ from those in  $R$  at no more than  $k$  registers, because flow  $f$  can only affect  $k$  registers. If  $k \ll m$ , the difference in the CDF of  $Z'$  compared to that of  $Z$  is small.
- Assume there are many flows in the stream and  $n_f$  is small compared to the total cardinalities of all other flows. As packets are being distributed to a register, it becomes harder and harder for the register's value to further increase, because it requires an much less likely (with geometrically decaying probability) hashed value to occur.

The analytical distribution of  $Z$  is also difficult to find. However, we can obtain an empirical distribution of  $Z$  directly from the register values in the pool at the end of the sketching process:

$$F_Z(i) \triangleq \Pr(Z \leq i) \approx \frac{\sum_{j=0}^{m-1} \mathbf{1}_{\{R[j] \leq i\}}}{m}, \quad 0 \leq i \leq r_{\max}. \quad (4.6)$$

In Figure 4.1 we show a sample distribution of  $Z$  generated using real trace data (the same data used for plotting Figure 2.1). From the figure we can see that the values of most of the registers are centered around 2 to 7.

Since  $Z'$  is the distribution of an arbitrary register value in  $R^{(1)}$ , including the  $k$  registers in  $R_f^{(2)}$ . We have  $Z_{f,0}^{(1)}, Z_{f,1}^{(1)}, \dots, Z_{f,k-1}^{(1)} \sim Z'$ . Since the  $k$  registers are randomly selected

from the pool with replacement, their independence is guaranteed by a good choice of the hash function  $G$ .

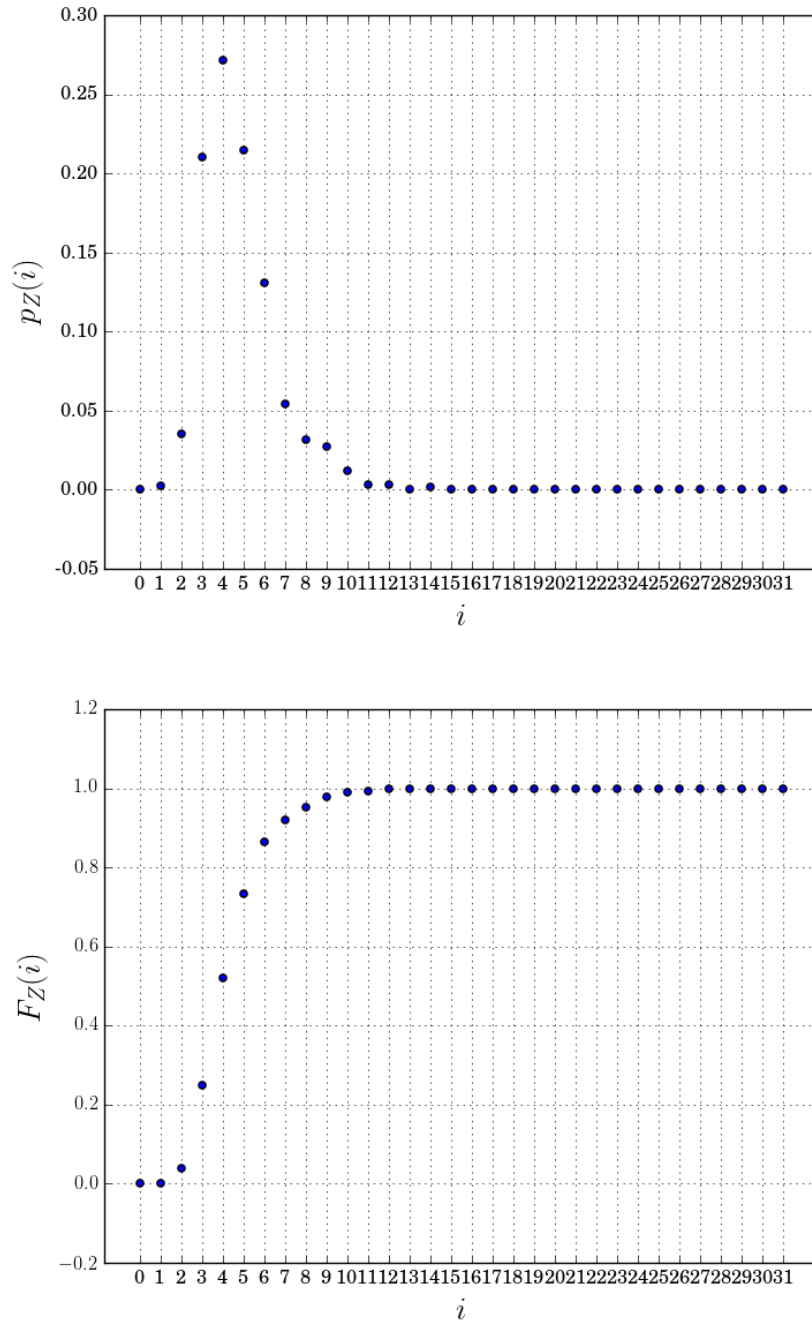


Figure 4.1: Sample distribution of  $Z$  generated using real trace data. Top: pmf of  $Z$ . Bottom: CDF of  $Z$ . In this case,  $m = 200,000$  and  $k = 512$ .

### 4.2.3 Distribution of $W_n$ and independence of $Z_{f,0}^{(2)}, \dots, Z_{f,k-1}^{(2)}$

To find the distribution of  $W_n$ , we first define a random variable  $X_n$ :

$X_n \triangleq$  the number of distinct elements distributed to an  
arbitrary register in  $R_f^{(2)}$ , conditioned on  $n_f = n$ .

It is easy to see  $X_n \sim \text{Binom}(n, \frac{1}{k})$ . The number of distinct elements distributed to each of the  $k$  registers in  $R_f^{(2)}$  has the same distribution as  $X_n$ , but they are not independent of each other (because they sum to  $n$ ). However, if  $n$  is large, we can approximate  $n$  by a Poisson random variable with mean  $n$ . Under this Poissonization approximation,  $X_n \sim \text{Poisson}(\frac{n}{k})$  and the number of distinct elements distributed to each of the  $k$  registers in  $R_f^{(2)}$  will then be independent (see Chapter 5.4 of [14] for this Poissonization trick).

As discussed in Section 2.1.1,  $W_n$  can be regarded as the maximum of  $X_n$  i.i.d. geometric random variables with parameter  $1/2$ . Given that  $X_n \sim \text{Poisson}(\frac{n}{k})$ , the CDF of  $W_n$  can be calculated as follows. For  $0 \leq i \leq r_{\max}$ ,

$$\begin{aligned} \Pr\{W_n \leq i\} &= \sum_{j=0}^n \Pr(X_n = j) \Pr\{\text{the maximum of } j \text{ i.i.d. Geo}(0.5) \text{ r.v.'s} \leq i\} \\ &= \sum_{j=0}^n \frac{(n/k)^j}{j!} e^{-n/k} \left(1 - \frac{1}{2^i}\right)^j \\ &= \sum_{j=0}^n \frac{\left(\frac{n}{k} \left(1 - \frac{1}{2^i}\right)\right)^j}{j!} e^{-n/k} \\ &= e^{-\frac{n}{k} \frac{1}{2^i}}. \end{aligned} \tag{4.7}$$

Under the Poissonization approximation,  $Z_{f,0}^{(2)}, Z_{f,1}^{(2)}, \dots, Z_{f,k-1}^{(2)}$  are independent and have the same distribution with  $W_n$ . We emphasize that this independence follows from the assumption that  $n$  is large.

### 4.2.4 Distribution of $Z_n$ and independence of $Z_{f,0}, \dots, Z_{f,k-1}$

Recall that, by construction,  $Z_n = \max\{Z', W_n\}$ . Therefore the CDF of  $Z_n$  is:

$$\Pr\{Z_n \leq i\} = \Pr\{\max\{Z', W_n\} \leq i\} = \Pr\{Z' \leq i\} \Pr\{W_n \leq i\} \approx F_Z(i) F_{W_n}(i), \quad 0 \leq i \leq r_{\max},$$

where  $F_Z(i)$  and  $F_{W_n}(i)$  are the CDF of  $Z$  and  $W_n$  respectively. We have used that  $Z'$  and  $W_n$  are independent of each other and the distribution of  $Z'$  can be approximated by that of  $Z$ . Hence, the pmf of  $Z_n$  can be approximated by

$$p_{Z_n}(i) = \begin{cases} F_Z(i)F_{W_n}(i) & \text{when } i = 0 \\ F_Z(i)F_{W_n}(i) - F_Z(i-1)F_{W_n}(i-1) & \text{when } 1 \leq i \leq r_{\max}, \end{cases} \quad (4.8)$$

where

$$F_{W_n}(i) = e^{-\frac{n}{k} \frac{1}{2^i}} \quad \text{with Poissonization of } n. \quad (4.9)$$

We have argued the independence of  $Z_{f,0}^{(1)}, \dots, Z_{f,k-1}^{(1)}$  and the independence of  $Z_{f,0}^{(2)}, \dots, Z_{f,k-1}^{(2)}$  under certain conditions. But  $Z_{f,i} = \max \{Z_{f,i}^{(1)}, Z_{f,i}^{(2)}\}$  and  $Z_{f,i}^{(1)}$  is independent of  $Z_{f,i}^{(2)}$  by construction. So  $Z_{f,0}, \dots, Z_{f,k-1}$  are also independent of each other under the same conditions. The consequence of this independence is that we can factor the  $k$ -variable likelihood function  $p_n(z_0, \dots, z_{k-1})$  into the product of  $k$  single-variable pmfs, based on the distribution of  $Z_n$ . That is, under this approximation, the log likelihood function in (4.2) can be re-written as

$$L_f(n) = \sum_{j=0}^{k-1} \ln p_{Z_n}(R_f[j]). \quad (4.10)$$

## 4.3 Implementation

It would be nice if we could obtain a closed-form analytical expression of the value of  $n$  that maximizes the log-likelihood function  $L_f(n)$ . Unfortunately it turns out to be difficult. In order to search for this value, we explore properties of  $L_f(n)$  that might be helpful.

### 4.3.1 Concavity of the log-likelihood function

In this subsection we show that  $L_f(n)$  has the decreasing increment property with respect to  $n$ . That is,  $L_f(n)$  is the restriction of a concave function to integer values. A by-product of this property is that  $L_f(n)$  must have a global maximum over the possible values of  $n$ . We will call this property “concave” or “concavity” for convenience henceforth.

Suppose for now that the possible values of  $n$  is a continuous range. First, since  $R_f[j] \in$

$\{0, \dots, r_{\max}\}$ , we have

$$L_f(n) = \sum_{j=0}^{k-1} \ln p_{Z_n}(R_f[j]) = \sum_{i=0}^{r_{\max}} c_i \ln p_{Z_n}(i), \quad (4.11)$$

for some non-negative integer constants  $c_0, c_1, \dots, c_{r_{\max}}$  such that  $\sum_{i=1}^{r_{\max}} c_i = k$  (i.e. the constants represent the empirical pmf of  $(R_f[0], \dots, R_f[k-1])$ ). Therefore a sufficient condition for  $L_f(n)$  to be concave in  $n$  is that  $\ln p_{Z_n}(i)$  is concave in  $n$  for each fixed value  $i \in \{0, \dots, r_{\max}\}$ . Recall the pmf of  $Z_n$  from (4.8) and (4.9). Let us denote  $a_i \triangleq F_Z(i)$  and  $b_i \triangleq e^{-\frac{1}{k2^i}}$ , so  $a_i, b_i \geq 0$ . Then we have

$$\ln p_{Z_n}(i) = \begin{cases} \ln a_i b_i^n & \text{if } i = 0 \\ \ln(a_i b_i^n - a_{i-1} b_{i-1}^n) & \text{if } 1 \leq i \leq r_{\max}. \end{cases}$$

When  $i = 0$ ,

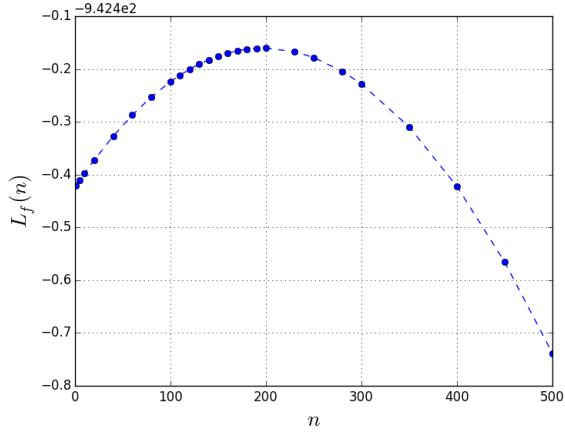
$$\ln p_n(i) = \ln a_i b_i^n = \ln a_i + n \ln b_i,$$

which is linear in  $n$  and thus concave in  $n$ . When  $1 \leq i \leq r_{\max}$ , we prove that  $\ln p_{Z_n}(i)$  is concave in  $n$  by showing that its second derivative with respect to  $n$  is non-positive:

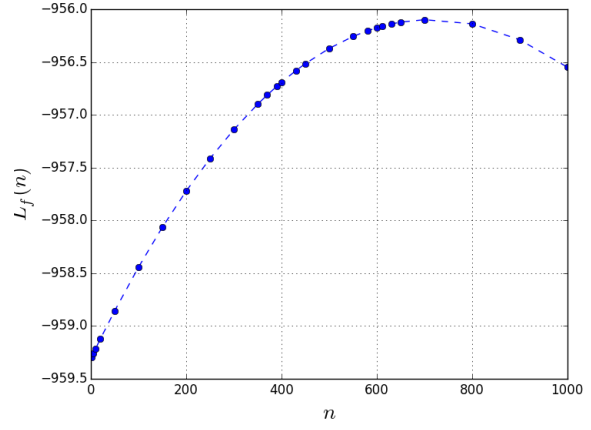
$$\begin{aligned} \ln p_n(i) &= \ln(a_i b_i^n - a_{i-1} b_{i-1}^n), \\ \frac{\partial \ln p_n(i)}{\partial n} &= \frac{a_i (\ln b_i) b_i^n - a_{i-1} (\ln b_{i-1}) b_{i-1}^n}{a_i b_i^n - a_{i-1} b_{i-1}^n}, \\ \frac{\partial^2 \ln p_n(i)}{\partial n^2} &= -\frac{a_{i-1} a_i b_{i-1}^n b_i^n (\ln b_{i-1} - \ln b_i)^2}{(a_i b_i^n - a_{i-1} b_{i-1}^n)^2} < 0. \end{aligned}$$

We conclude that  $L_f(n)$  is concave in  $n$ .

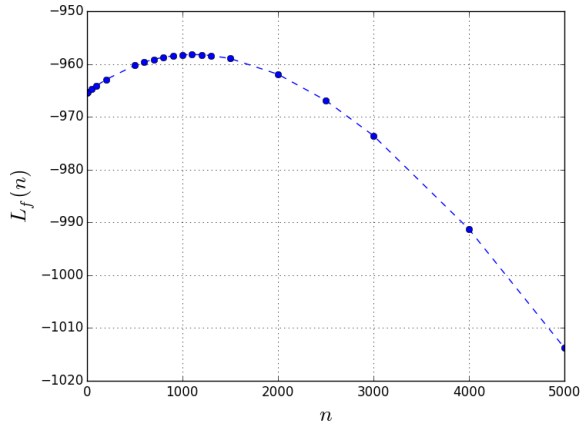
In Figure 4.2, we show the plot of  $L_f(n)$  as a function of  $n$  for four selected flows whose cardinalities are 150, 611, 1000 and 31536, respectively. The plots are generated based on the same trace data used for Figure 2.1 and Figure 4.1. The concavity of  $L_f(n)$  is obvious. We can also observe that in each plot, the value of  $n$  at which the curve peaks is close to the actual flow cardinality  $n_f$  (though not exact the same due to randomness and noise).



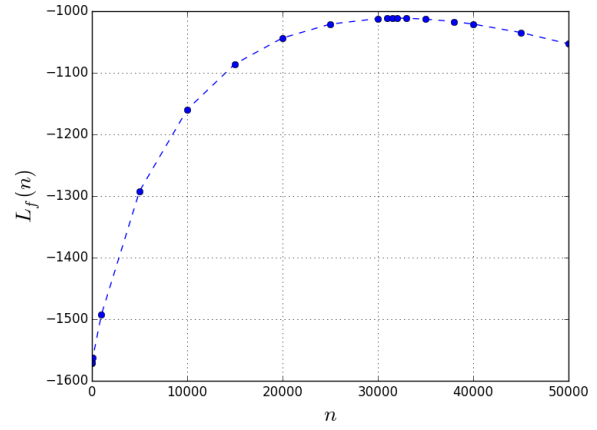
(a)  $n_f = 150$



(b)  $n_f = 611$



(c)  $n_f = 1000$



(d)  $n_f = 31536$

Figure 4.2: Four sample plots of log likelihood function  $L_f(n)$  vs.  $n$ .

### 4.3.2 Bi-section search implementation

We have shown that  $L_f(n)$  is concave in  $n$ . Since  $n$  can only take integer values in  $\{1, 2, \dots, n_{\max}\}$ , we can do a bi-section search to find the optimal  $n$ .

**Bi-section search implementation of  $\hat{n}_{f,ML}$ :**

1.  $lb \leftarrow 1, ub \leftarrow n_{\max}$ .
2. While  $ub - lb > 1$ :
  - $mid_1 \leftarrow \lfloor \frac{lb+ub}{2} \rfloor, mid_2 \leftarrow mid_1 + 1$
  - if  $L_f(mid_1) \geq L_f(mid_2)$ :
    - $lb \leftarrow lb, ub \leftarrow mid_1$
  - else:
    - $lb \leftarrow mid_2, ub \leftarrow ub$
3. If  $L_f(lb) \geq L_f(ub)$ , return  $lb$ ; otherwise return  $ub$ .

The pseudo-code above should be self-explanatory. In actual implementation, it is found that the algorithm reaches the optimal  $n$  at a faster rate if we replace  $mid_1 = \lfloor \frac{lb+ub}{2} \rfloor$  with  $mid_1 = \sqrt{lb \cdot ub}$ , which is equivalent to  $\log mid_1 = \frac{1}{2}(\log lb + \log ub)$ , i.e. a bi-section search on log scale. This faster rate is because, as we already mentioned, most of the flows have small cardinalities. If  $n_{\max} = 10^6$ , then at the first while loop iteration, with  $mid_1 = \lfloor \frac{lb+ub}{2} \rfloor$ , we have  $mid_1 = 5 \times 10^5$ ; but with  $mid_1 = \sqrt{lb \cdot ub}$  we have  $mid_1 = 1000$  — the latter search is quicker if the optimal  $n$  is actually small.

With bi-section search, the number of searches required for the estimation of one flow is upper bounded by  $\log n_{\max}$ . This is not ideal compared to the  $vLL_{\theta}$  estimator which only needs one search/estimation. However, since the estimation process is off-line, the worst-case time complexity of  $\log n_{\max}$  for searching is still acceptable.

### 4.3.3 Summary of the maximum-likelihood estimator

We wrap up the ideas in the previous sections of this chapter and name this estimator  $vLL$ -MLE, meaning that it uses *maximum likelihood estimation* based on *virtual LogLog sketching*.



The estimator is summarized as follows. To estimate the cardinality of flow  $f$ :

1. Find the values of the  $k$  registers in  $R_f$ :  $R_f[0], \dots, R_f[k-1]$ .
2. From the register values in  $R$ , find the CDF of  $Z$  by (4.6).
3. Do a bi-section search on integer value  $n$  in the range  $\{1, 2, \dots, n_{\max}\}$  to maximize the log likelihood function  $L_f(n)$ , output this value of  $n$  as the estimate.
  - For a given  $n$ ,  $L_f(n)$  can be evaluated by (4.8), (4.9) and (4.10).

## 4.4 Aside: MLE for single-flow estimation

As an aside, we can also use the vLL-MLE estimator for single-flow estimation, with some minor changes. Recall that in single-flow estimation, the goal is to estimate the flow's cardinality  $n$  from its  $k$  register values. This is similar to the per-flow estimation case, where the goal is to estimate a given flow  $f$ 's cardinality  $n_f$  from its  $k$  register values  $R_f[0], \dots, R_f[k-1]$ , except that in the single-flow case, there is no background noise in the register values. That is, the vLL-MLE estimator for per-flow estimation can be used for single-flow estimation by simply letting  $Z = 0$ .

In this case, we do not need the assumptions related to  $Z$  anymore. But the Poissonization approximation of  $n$  (i.e. assuming  $n$  is large) is still necessary.

To evaluate the performance of the MLE on single-flow estimation, we run experiments to obtain empirical values of its relative standard error for selected values of  $k$  and compare it with the HLL estimator. The experiments are performed with  $n = 10^7$ . The results are summarized in Table 4.1.

Table 4.1: Compare the accuracy (relative standard error) of the MLE and HLL estimators with selected values of  $k$ .

Estimator	$k = 512$	$k = 1024$	$k = 2048$
MLE	0.04609	0.03127	0.02220
HLL	0.04585	0.03127	0.02216

The experiment results show that the MLE's performance is very close to that of HLL, but not any better. It reinforces the claim that the HLL estimator is near-optimal for single-flow estimation.

## 4.5 Experimental performance evaluation

In this section we evaluate the performance of the vLL-MLE estimator and compare it with that of the  $vLL_\theta$  estimator introduced in Chapter 3. Since we have shown that  $vLL_{-1}$  (i.e. vHLL) is the best within the family of the  $vLL_\theta$  estimators, we compare vLL-MLE with  $vLL_{-1}$  in particular. Results for each estimator are generated by running experiments on the same 100 simulated trace files.

Figure 4.3 shows the estimation results of both estimators by directly plotting the estimated cardinalities vs. the corresponding actual cardinalities. Each point in the graphs represents one flow, with its x-coordinate value being the actual cardinality of the flow and its y-coordinate value being the estimated cardinality. The more clustered the points are to the equality line  $y = x$ , the more accurate the estimator is. We can see that the two estimators have comparable performances.

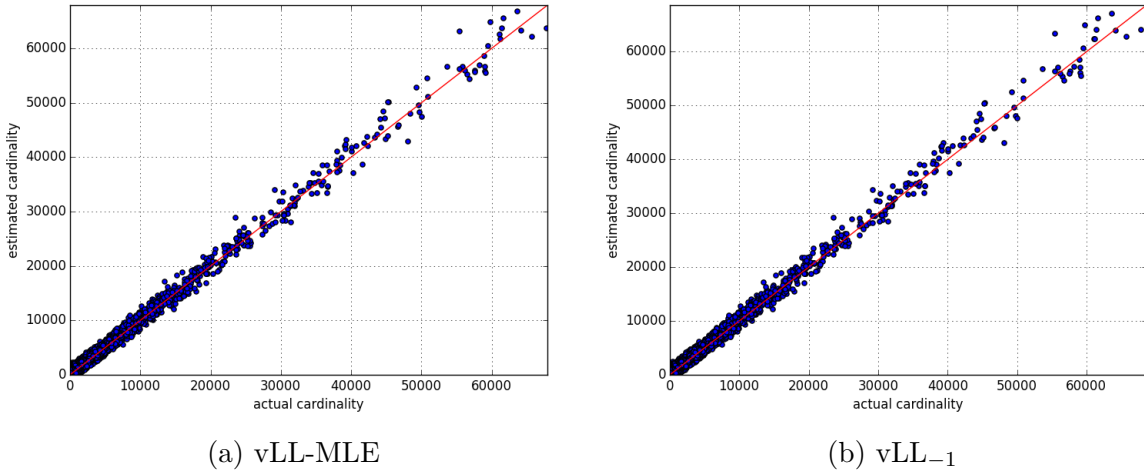


Figure 4.3: Plot of estimated cardinalities vs. actual cardinalities for the vLL-MLE and  $vLL_{-1}$  estimators.

In Figure 4.4 and Figure 4.5, we respectively plot the relative bias (defined as  $\mathbb{E} \left( \frac{\hat{n}_f}{n_f} \right) - 1$ ) and relative standard error (defined as  $\frac{\sqrt{\text{Var}(\hat{n}_f)}}{n_f}$ ) vs. the actual flow cardinality for both estimators. Since there are not many flows for some cardinalities (especial the large cardinalities), we divide the horizontal axis into bins of width 1000 for cardinalities  $\leq 10000$  and of width 5000 for cardinalities  $> 10000$ . In each bin, we calculate the empirical relative bias and relative standard error of the data and interpolate the values for each bin on the graph to form the plots. Again we see that the two estimators have very similar performances.

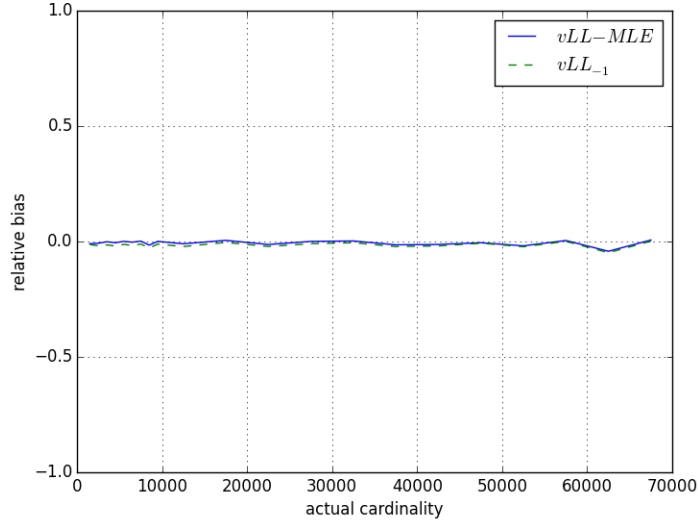


Figure 4.4: Simulated results on the relative bias of the  $vLL-MLE$  and  $vLL_{-1}$  estimators.

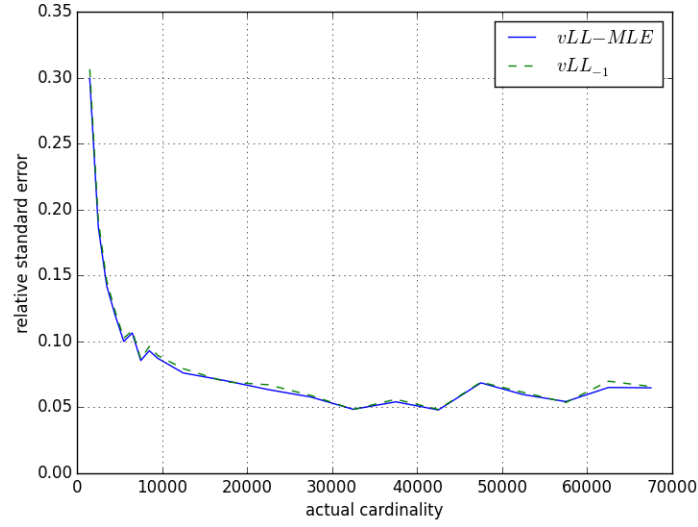


Figure 4.5: Simulated results on the relative standard error of the  $vLL-MLE$  and  $vLL_{-1}$  estimators.

In both figures, the curves are plotted for cardinalities larger than 1000 only. This is for better graph layouts: we found from the experiment results that estimates for flows with cardinalities less than 1000 are very inaccurate. Accurate estimation for small flows is difficult because of the noise caused by large flows in their register values.

Figure 4.4 shows that both estimators are approximately unbiased for large flows (with

cardinalities  $> 1000$ ).

Figure 4.5 shows that in general, for both estimators, the estimation accuracy improves as the cardinality gets larger. We do observe a slight curving up at the high end of the graph though, meaning that the accuracy stops improving once the actual flow cardinality reaches a certain level (probably  $> 40000$ ).

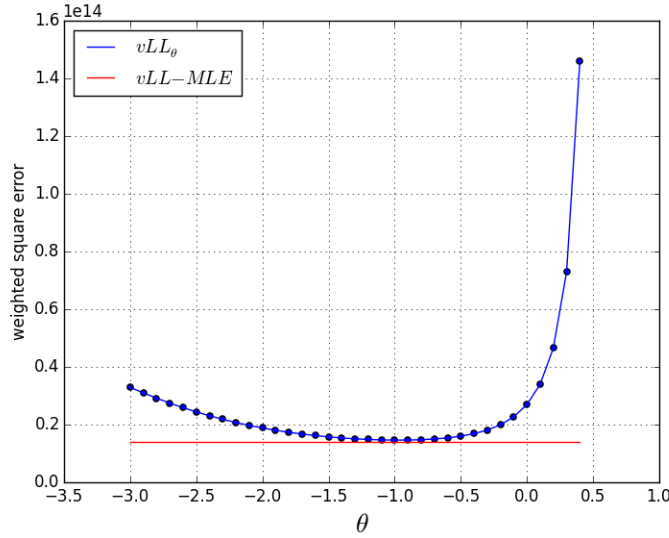


Figure 4.6: Compare the weighted square error of the vLL-MLE estimator and the vLL $_{\theta}$  estimators.

Finally, in Figure 4.6 we compare the weighted square error of the vLL-MLE estimator with that of the vLL $_{\theta}$  estimators. Experiment results show that the vLL-MLE estimator outperforms the vLL $_{\theta}$  estimator for all values of  $\theta$ ; compared to the vLL $_{-1}$  estimator, the vLL-MLE estimator has a slight improvement of about 3.5%. We conclude that the two estimators, vLL-MLE and vLL $_{-1}$  (i.e. vHLL), have comparable performances.

## Chapter 5

# CONCLUSIONS AND FUTURE WORK

In this thesis we explored two new perspectives on the estimation process of the virtual LogLog algorithm [3] for per-flow cardinality estimation:

- We showed how the existing vHLL estimator of [3] for per-flow estimation can be generalized by introducing a parameter  $\theta$ , in a similar way in which the HLL estimator for single-flow estimation can be generalized by  $\theta$ .
- We proposed the vLL-MLE estimator, an alternative approach to the per-flow estimation problem.

In both cases we provided empirical evidence to show the near-optimality of the vHLL estimator for per-flow estimation. This result is analogous to the near-optimality of the HLL estimator for single-flow estimation [7].

Results of this thesis are mostly based on simulated experiments. One possible future work is the analysis of theoretical bounds for per-flow cardinality estimation. For example, for a given amount of memory, flow cardinality distribution and number of flows, what is the best possible level of accuracy that can be achieved for per-flow estimation? There has been much research on this direction for single-flow estimation. For example, [2] gives an asymptotic lower bound on the relative standard error of single-flow estimators based on order statistics. It would be useful to obtain similar results for per-flow estimation.

Another possible direction of future work is on efficient algorithms that identify heavy-hitters (i.e. large flows). With the per-flow estimators discussed in this thesis, one is able to estimate the cardinality of a flow *given* the flow's ID. However, in many applications, the goal is to *identify* large flows, in which case we are not given the IDs of these flows beforehand. One possible approach to this problem is to store all the distinct flow IDs using a separate block of memory and then check each flow one by one; the work in [16] discusses how all the distinct flow IDs can be stored in main memory. However, with this approach some memory is wasted on flows that are actually small (which do not matter at all). Algorithms that

identify large flows and estimate their cardinalities directly and more efficiently are useful in this context.

# REFERENCES

- [1] F. Giroire, “Order statistics and estimating cardinalities of massive data sets,” *Discrete Applied Mathematics*, vol. 157, pp. 406–427, Jan 2009.
- [2] P. Chassaing and L. Gerin, “Efficient estimation of the cardinality of large data sets,” 2011, arXiv:math/0701347v3.
- [3] Q. Xiao, S. Chen, M. Chen, and Y. Ling, “Hyper-compact virtual estimators for big network data based on register sharing,” in *Proceedings of the ACM SIGMETRICS 2015*, Portland, Oregon, USA, 2015, pp. 417–428.
- [4] C. Estan, G. Varshese, and M. Fisk, “Bitmap algorithms for counting active flows on high speed link,” *IEEE/ACM Transactions on Networking*, vol. 14, pp. 923–937, Oct 2006.
- [5] P. Flajolet and G. N. Martin, “Probabilistic counting algorithms for data base algorithms,” *Journal of Computer and System Sciences*, vol. 31, no. 2, pp. 182–209, Sep 1985.
- [6] S. Heule, M. Nunkesser, and A. Hall, “Hyperloglog in practice: Algorithmic engineering of a state of the art cardinality estimation algorithm,” in *Proceedings of the EDBT 2013 Conference*, Genoa, Italy, Mar 2013.
- [7] P. Flajolet, E. Fusy, O. Gandouet, and F. Meunier, “Hyperloglog: The analysis of a near-optimal cardinality estimation algorithm,” in *Proceedings of AOFA ’07*, 2007, pp. 127–146.
- [8] M. Durand and P. Flajolet, “Loglog counting of large cardinalities,” in *European Symposium on Algorithms*, 2003, pp. 605–617.
- [9] P. B. Gibbons, “Distinct-values estimation over data streams,” in *Data Streams Management: Processing High-Speed Data Streams*, 2007. [Online]. Available: <http://www.pittsburgh.intel-research.net/people/gibbons/>
- [10] A. Metwally, D. Agrawal, and A. E. Abbadi, “Why go logarithmic if we can go linear?: Towards effective distinct counting of search traffic,” in *Proceedings of the 11th International Conference on Extending Database Technology: Advances in Database Technology*, ser. EDBT ’08. New York, NY, USA: ACM, 2008, pp. 618–629.

- [11] P. Clifford and I. A. Cosma, “A statistical analysis of probabilistic counting algorithms,” *Scandinavian Journal of Statistics*, vol. 39, no. 1, pp. 1–14, Jun 2010.
- [12] “The CAIDA UCSD anonymized internet traces 2013 on Jan 17.” [http://www.caida.org/data/passive/passive\\_2013\\_dataset.xml](http://www.caida.org/data/passive/passive_2013_dataset.xml).
- [13] Z. Mo, Y. Qiao, S. Chen, and T. Li, “Highly compact virtual maximum likelihood sketches for counting big network data,” in *52 Annual Allerton Conference*, Allerton House, UIUC, Illinois, USA, Oct 2014.
- [14] M. Mitzenmacher and E. Upfal, *Probability and Computing: Randomized Algorithms and Probabilistic Analysis*, 1st ed. Cambridge, UK: Cambridge University Press, 2005.
- [15] L. A. Adamic and B. A. Huberman, *Glottometrics*, pp. 143–150, 2002.
- [16] M. Yoon, T. Li, S. Chen, and J.-K. Peir, “Fit a compact spread estimator in small high-speed memory,” *IEEE/ACM Trans. Netw.*, vol. 19, no. 5, pp. 1253–1264, Oct 2011.