

---

# Composable Sketches for Functions of Frequencies: Beyond the Worst Case

---

Edith Cohen<sup>1,2</sup> Ofir Geri<sup>3,4</sup> Rasmus Pagh<sup>1,5,6</sup>

## Abstract

Recently there has been increased interest in using machine learning techniques to improve classical algorithms. In this paper we study when it is possible to construct compact, composable sketches for weighted sampling and statistics estimation according to functions of data frequencies. Such structures are now central components of large-scale data analytics and machine learning pipelines. However, many common functions, such as thresholds and  $p$ th frequency moments with  $p > 2$ , are known to require polynomial size sketches in the worst case. We explore performance beyond the worst case under two different types of assumptions. The first is having access to noisy *advice* on item frequencies. This continues the line of work of Hsu et al. (ICLR 2019), who assume predictions are provided by a machine learning model. The second is providing guaranteed performance on a restricted class of input frequency distributions that are better aligned with what is observed in practice. This extends the work on heavy hitters under Zipfian distributions in a seminal paper of Charikar et al. (ICALP 2002). Surprisingly, we show analytically and empirically that “in practice” small polylogarithmic-size sketches provide accuracy for “hard” functions.

## 1. Introduction

Composable sketches, also known as mergeable summaries (Agarwal et al., 2013), are data structures that support summarizing large amounts of distributed or streamed

---

<sup>1</sup>Google Research, Mountain View, CA <sup>2</sup>Department of Computer Science, Tel Aviv University, Israel <sup>3</sup>Stanford University, CA, USA <sup>4</sup>Most of this work was done while interning at Google Research. <sup>5</sup>BARC, Denmark <sup>6</sup>IT University of Copenhagen, Denmark. Correspondence to: Edith Cohen <edith@cohenwang.com>, Ofir Geri <ofirgeri@cs.stanford.edu>, Rasmus Pagh <pagh@itu.dk>.

data with small computational resources (time, communication, and space). Such sketches support processing additional data elements and merging sketches of multiple datasets to obtain a sketch of the union of the datasets. This design is suitable for working with streaming data (by processing elements as they arrive) and distributed datasets, and allows us to parallelize computations over massive datasets. Sketches are now a central part of managing large-scale data, with application areas as varied as federated learning (McMahan et al., 2017) and statistics collection at network switches (Liu et al., 2016; 2019).

The datasets we consider consist of *elements* that are key-value pairs  $(x, v)$  where  $v \geq 0$ . The frequency  $w_x$  of a key  $x$  is defined as the sum of the values of elements with that key. When the value of all elements is 1, the frequency is simply the number of occurrences of a key in the dataset. Examples of such datasets include search queries, network traffic, user interactions, or training data from many devices. These datasets are typically distributed or streamed.

Given a dataset of this form, one is often interested in computing statistics that depend on the frequencies of keys. For example, the statistics of interest can be the number of keys with frequency greater than some constant (threshold functions), or the second frequency moment ( $\sum_x w_x^2$ ), which can be used to estimate the skew of the data. Generally, we are interested in statistics of the form

$$\sum_x L_x f(w_x) \quad (1)$$

where  $f$  is some function applied to the frequencies of the keys and the coefficients  $L_x$  are provided (for example as a function of the features of the key  $x$ ). An important special case, popularized in the seminal work of (Alon et al., 1999), is computing the  $f$ -statistics of the data:  $\|f(\mathbf{w})\|_1 = \sum_x f(w_x)$ .

One way to compute statistics of the form (1) is to compute a random sample of keys, and then use the sample to compute estimates for the statistics. In order to compute low-error estimates, the sampling has to be weighted in a way that depends on the target function  $f$ : each key  $x$  is weighted by  $f(w_x)$ . Since the problem of computing a weighted sample is more general than computing  $f$ -statistics, our focus in this work will be on composable sketches for weighted sampling according to different functions of frequencies.

The tasks of sampling or statistics computation can always be performed by first computing a table of key and frequency pairs  $(x, w_x)$ . But this aggregation requires a data structure of size (and in turn, communication or space) that grows linearly with the number of keys whereas ideally we want the size to grow at most polylogarithmically. With such small sketches we can only hope for approximate results and generally we see a trade-off between sketch size, which determines the storage or communication needs of the computation, and accuracy.

When estimating statistics from samples, the accuracy depends on the sample size and on how much the sampling probabilities “suit” the statistics we are estimating. In order to minimize the error, the sampling probability of each key  $x$  should be (roughly) proportional to  $f(w_x)$ . This leads to a natural and extensively-studied question: for which functions  $f$  can we design efficient sampling sketches?

The literature and practice are ripe with surprising successes for sketching, including small (polylogarithmic size) sketch structures for estimating the number of distinct elements (Flajolet & Martin, 1985; Flajolet et al., 2007) ( $f(w) = I_{w>0}$ ), frequency moments ( $f(w) = w^p$ ) for  $p \in [0, 2]$  (Alon et al., 1999; Indyk, 2001), and computing  $\ell_p$  heavy hitters (for  $p \leq 2$ , where an  $\ell_p$   $\varepsilon$ -heavy hitter is a key  $x$  with  $w_x^p \geq \varepsilon \|w\|_p^p$ ) (Misra & Gries, 1982; Charikar et al., 2002; Manku & Motwani, 2002; Cormode & Muthukrishnan, 2005; Metwally et al., 2005). Here  $I_\sigma$  is the indicator function that is 1 if the predicate  $\sigma$  is true, and 0 otherwise. A variety of methods now support sampling via small sketches for rich classes of functions of frequencies (Cohen, 2018; McGregor et al., 2016; Jayaram & Woodruff, 2018; Cohen & Geri, 2019), including the moments  $f(w) = w^p$  for  $p \in [0, 2]$  and the family of concave sublinear functions.

The flip side is that we know of lower bounds that limit the performance of sketches using small space for some fundamental tasks (Alon et al., 1999). A full characterization of functions  $f$  for which  $f$ -statistics can be estimated using polylogarithmic-size sketches was provided in (Braverman & Ostrovsky, 2010). Examples of “hard” functions are thresholds  $f(w) = I_{w>T}$  (counting the number of keys with frequency above a specified threshold value  $T$ ), threshold weights  $f(w) = wI_{w>T}$ , and high frequency moments  $f(w) = w^p$  with  $p > 2$ . Estimating the  $p$ th frequency moment ( $\sum_x w_x^p$ ) for  $p > 2$  is known to require space  $\Omega(n^{1-2/p})$  (Alon et al., 1999; Li & Woodruff, 2013), where  $n$  is the number of keys. These particular functions are important for downstream tasks: Threshold aggregates characterize the distribution of frequencies, and high moment estimation is used in the method of moments, graph applications (Eden et al., 2019), and for estimating the cardinality of multi-way self-joins (Alon et al., 2002) (a  $p$ th moment is used for estimating a  $p$ -way join).

**Beyond the worst case.** Much of the discussion of sketching classified functions into “easy” and “hard”. For example, there are known efficient methods for sampling according to  $f(w) = w^p$  for  $p \leq 2$ , while for  $p > 2$ , even the easier task of computing the  $p$ th moment is known to require polynomial space. However, the hard data distributions used to establish lower bounds for some functions of frequency are arguably not very realistic. Real data tends to follow nice distributions and is often (at least somewhat) predictable. We study sketching where additional assumptions allow us to circumvent these lower bounds while still providing theoretical guarantees on the quality of the estimates. We consider two distinct ways of going beyond the worst case: 1) access to *advice models*, and 2) making natural assumptions on the frequency distribution of the dataset.

For the sampling sketches described in this paper, we use a notion of *overhead* to capture the discrepancy between the sampling probabilities used in the sketch and the “ideal” sampling probabilities of weighted sampling according to a target function of frequency  $f$ . An immensely powerful property of using sampling to estimate statistics of the form (1) is that the overhead translates into a multiplicative increase in sample/sketch size, without compromising the accuracy of the results (with respect to what an ideal “benchmark” weighted sample provides). This property was used in different contexts in prior work, e.g., (Frieze et al., 2004; Cohen et al., 2009), and we show that it can be harnessed for our purposes as well. For the task of estimating  $f$ -statistics, we use a tailored definition of overhead, that is smaller than the overhead for the more general statistics (1).

**Advice model.** The advice model for sketching was recently proposed and studied by Hsu et al. (2019). The advice takes the form of an oracle that is able to identify whether a given key is a heavy hitter. Such advice can be generated, for example, by a machine learning model trained on past data. The use of the “predictability” of data to improve performance was also demonstrated in (Kraska et al., 2018; Indyk et al., 2019). A similar heavy hitter oracle was used in (Jiang et al., 2020) to study additional problems in the streaming setting. For high frequency moments, they obtained sketch size  $O(n^{1/2-1/p})$ , a quadratic improvement over the worst-case lower bound.

Here we propose a sketch for *sampling by advice*. We assume an advice oracle that returns a noisy prediction of the frequency of each key. This type of advice oracle was used in the experimental section of (Hsu et al., 2019) in order to detect heavy hitters. We show that when the predicted  $f(w_x)$  for keys with above-average contributions  $f(w_x)$  is approximately accurate within a factor  $C$ , our sample has overhead  $O(C)$ . That is, the uncertainty in the advice translates to a factor  $O(C)$  increase in the sketch size but does *not* impact the accuracy.

**Frequencies-functions combinations.** Typically, one designs sketch structures to provide guarantees for a certain function  $f$  and any set of input frequencies  $\mathbf{w}$ . The performance of a sketch structure is then analyzed for a *worst-case* frequency distribution. The analysis of the advice model also assumes worst-case distributions (with the benefits that comes from the advice). We depart from this and study sketch performance for a *combination*  $(F, W, h)$  of a family of functions  $F$ , a family  $W$  of frequency distributions, and an overhead factor  $h$ . Specifically, we seek sampling sketches that produce weighted samples with overhead at most  $h$  with respect to  $f(\mathbf{w})$  for *every* function  $f \in F$  and frequency distribution  $\mathbf{w} \in W$ . By limiting the set  $W$  of input frequency distributions we are able to provide performance guarantees for a wider set  $F$  of functions of frequency, including functions that are worst-case hard. We particularly seek combinations with frequency distributions  $W$  that are typical in practice. Another powerful property of the combination formulation is that it provides multi-objective guarantees with respect to a multiple functions of frequency  $F$  using the same sketch (Cohen, 2015; 2018; Liu et al., 2016).

The performance of sketch structures on “natural” distributions was previously considered in a seminal paper by Charikar et al. (2002). The paper introduced the *Count Sketch* structure for heavy hitter detection, where an  $\ell_2$   $\varepsilon$ -heavy hitter is defined to be a key with  $w_x^2 \geq \varepsilon \|\mathbf{w}\|_2^2$ . They also show that for Zipf-distributed data with parameter  $\alpha \geq 1/2$ , a count sketch of size  $O(k)$  can be used to find the  $k$  heaviest keys (a worst-case hard problem) and that an  $\ell_1$  sample can only identify the heaviest keys for Zipf parameter  $\alpha \geq 1$ .

We significantly extend these insights to a wider family of frequency distributions and to a surprisingly broad class of functions of frequencies. In particular we show that all high moments ( $p > 2$ ) are “easy” as long as the frequency distribution has an  $\ell_1$  or  $\ell_2$   $\varepsilon$ -heavy hitter. In this case, an  $\ell_1$  or  $\ell_2$  sample with overhead  $1/\varepsilon$  can be used to estimate all high moments. We also show that in a sense this characterization is tight in that if we allow all frequencies, we meet the known lower bounds. It is very common for data sets in practice to have a most frequent key that is an  $\ell_1$  or  $\ell_2$   $\varepsilon$ -heavy hitter. This holds in particular for Zipf or approximate Zipf distributions.

Moreover, we show that Zipf frequency distributions have small *universal* sketches that apply to *any* monotone function of frequency (including thresholds and high moments). Zipf frequencies were previously considered in the advice model (Aamand et al., 2019). Interestingly, we show that for these distribution a single small sketch is effective with all monotone functions of frequency, even without advice. In these cases, universal sampling is achieved with off-the-

shelf polylogarithmic-size sketches such as  $\ell_p$  samples for  $p \leq 2$  and multi-objective concave-sublinear samples (Cohen et al., 2012; Cohen, 2018; McGregor et al., 2016; Jayaram & Woodruff, 2018).

**Empirical study.** We complement our analysis with an empirical study on multiple real-world datasets including datasets studied in prior work on advice models (Pass et al., 2006; CAIDA, 2016; Paranjape et al., 2017; Maciá-Fernández et al., 2018). (Additional discussion of the datasets appears in the Appendix A.) We apply sampling by advice, with advice based on models from prior work or direct use of frequencies from past data. We then estimate high moments from the samples. We observe that sampling-by-advice was effective on these datasets, yielding low error with small sample size. We also observed however that  $\ell_2$  and  $\ell_1$  samplers were surprisingly accurate on these tasks, with  $\ell_2$  samplers generally outperforming sampling by advice. This surprisingly good performance of these simple sampling schemes is suggested from our analysis.

We compute the overhead factors for some off-the-shelf sampling sketches on multiple real-world datasets with the objectives of  $\ell_p$  sampling ( $p > 2$ ) and universal sampling. We find these factors to be surprisingly small. For example, the measured overhead of using  $\ell_2$  sampling for the objective of  $\ell_p$  sampling is in the range  $[1.2, 18]$ . For universal sampling, the observed overhead is lower with  $\ell_1$  and with multi-objective concave sublinear samples than with  $\ell_2$  sampling and is in  $[93, 800]$ , comparing very favorably with the alternative of computing a full table. Finally, we use sketches to estimate the distribution of rank versus frequency, which is an important tool for optimizing performance across application domains (for network flows, files, jobs, or search queries). We find that  $\ell_1$  samples provide quality estimates, which is explained by our analytical results.

## 2. Preliminaries

We consider datasets where each data element is a (key, value) pair. The keys belong to a universe denoted by  $\mathcal{X}$  (e.g., the set of possible users or words), and each key may appear in more than one element. The values are positive, and for each key  $x \in \mathcal{X}$ , we define its *frequency*  $w_x \geq 0$  to be the sum of values of all elements with key  $x$ . The data elements may appear as a stream or be stored in a distributed manner. We denote the number of active keys (keys with frequency greater than 0) by  $n$ .

We are interested in sketches that produce a weighted sample of keys according to some function  $f$  of their frequencies, where roughly key  $x$  is sampled with probability proportional to  $f(w_x)$ . We use  $f(\mathbf{w})$  as a shorthand for the vector of all values  $f(w_x)$  (in any fixed order).

**Estimates from a sample.** We work with sampling schemes that produce a random subset  $S \subseteq \mathcal{X}$  of the keys in the dataset. For each key  $x \in S$  we have its frequency  $w_x$  and can compute its probability  $p_x$  to be sampled. From such a sample, we can compute for each key  $x$  the *inverse probability estimate* (Horvitz & Thompson, 1952) of  $f(w_x)$  defined as

$$\widehat{f(w_x)} = \begin{cases} \frac{f(w_x)}{p_x} & \text{if } x \in S \\ 0 & \text{if } x \notin S \end{cases}.$$

These *unbiased* per-key estimates can be summed to obtain unbiased estimates of the  $f$ -statistics of a domain  $H \subset \mathcal{X}$ :

$$\sum_{x \in H} \widehat{f(w_x)} := \sum_{x \in H} \widehat{f(w_x)} = \sum_{x \in S} \widehat{f(w_x)}.$$

The last equality follows because  $\widehat{f(w_x)} = 0$  for keys not in the sample. We can similarly estimate other statistics that are linear in  $f(w_x)$ , e.g.,  $\sum_{x \in \mathcal{X}} L_x f(w_x)$  (for coefficients  $L_x$ ).

**Benchmark variance bounds.** We measure performance with respect to that of a “benchmark” weighted sampling scheme where each key  $x$  is sampled with probability proportional to  $f(w_x)$ . Recall that for “hard” functions  $f$  these schemes can not be implemented with small sketches. These benchmark schemes include (i) *probability proportional to size* (pps) with replacement where we have  $k$  independent draws where key  $x$  is selected with probability  $f(w_x)/\|f(\mathbf{w})\|_1$ , (ii) pps without replacement (ppswor (Rosén, 1972; 1997; Cohen & Kaplan, 2008)), or (iii) priority sampling (Ohlsson, 1998; Duffield et al., 2007). With these schemes, the variance for key  $x$  is upper bounded by<sup>1</sup>

$$\text{Var}[\widehat{f(w_x)}] \leq \frac{1}{k} f(w_x) \|f(\mathbf{w})\|_1 \quad (2)$$

Consequently, the variance of the sum estimator for the  $f$ -statistics of a domain  $H$  is bounded by (due to nonpositive covariance shown in earlier works):

$$\text{Var} \left[ \sum_{x \in H} \widehat{f(w_x)} \right] \leq \frac{1}{k} \sum_{x \in H} f(w_x) \|f(\mathbf{w})\|_1. \quad (3)$$

<sup>1</sup>With ppswor and priority sampling,  $k - 2$  instead of  $k$ . We use this upper bound as the benchmark, since when the sampling probability  $p_x = f(w_x)/\|f(\mathbf{w})\|_1$  approaches 1 (one key dominates the data), the variance  $(f(w_x))^2 \left(\frac{1}{p_x} - 1\right)$  approaches 0, and we cannot approximate it multiplicatively with a sampling probability that multiplicatively approximates  $p_x$ . However, when there is not just one key that dominates the data, for example,  $p_x \leq 1/2$  for all  $x$ , or more generally,  $p_x \leq 1 - \frac{1}{c}$ , we get that  $\text{Var}[\widehat{f(w_x)}] \geq (f(w_x))^2 \frac{1}{cp_x}$ , so the bound on the variance is almost tight.

The variance on the estimate of  $\|f(\mathbf{w})\|_1$  is bounded by

$$\text{Var}[\|\widehat{f(\mathbf{w})}\|_1] \leq \frac{1}{k} \|f(\mathbf{w})\|_1^2.$$

With these “benchmark” schemes, if we wish to get multiplicative error bound (normalized root mean squared error) of  $\varepsilon$  for estimating  $\|f(\mathbf{w})\|_1$  we need sample size  $k = O(\varepsilon^{-2})$ . We note that the estimates are also concentrated in the Chernoff sense (Duffield et al., 2007; Cohen, 2015).

We refer to the probability vector

$$p_x := \frac{f(w_x)}{\|f(\mathbf{w})\|_1}$$

as *pps sampling probabilities* for  $f(\mathbf{w})$ . When  $f(w) = w^p$  (for  $p > 0$ ) we refer to sampling with the respective pps probabilities as  $\ell_p$  sampling.

**Emulating a weighted sample.** Let  $\mathbf{p}$  be the base pps probabilities for  $f(\mathbf{w})$ . When we use a weighted sampling scheme with weights  $\mathbf{q} \neq \mathbf{p}$  then the variance bound (3) does not apply. We will say that weighted sampling according to  $\mathbf{q}$  *emulates* weighted sampling according to  $\mathbf{p}$  with *overhead*  $h$  if for all  $k$  and for all  $H$ , a sample of size  $kh$  provides the variance bound (3) (and the respective concentration bounds).

**Lemma 2.1.** *The overhead of emulating weighted sample according to  $\mathbf{p}$  using weighted sampling according to  $\mathbf{q}$  is at most*

$$h(\mathbf{p}, \mathbf{q}) := \max_x \frac{p_x}{q_x}.$$

*Proof.* We first bound the variance of  $\widehat{f(w_x)}$  for a key  $x$  with weighted sampling by  $\mathbf{q}$ . Consider a weighted sample of size  $k$  according to base probabilities  $q_x$ . Then for all  $x$ ,

$$\begin{aligned} \text{Var}[\widehat{f(w_x)}] &\leq \frac{1}{k} (f(w_x))^2 \left(\frac{1}{q_x} - 1\right) \leq \frac{1}{k} (f(w_x))^2 \frac{1}{p_x} \frac{p_x}{q_x} \\ &\leq \frac{1}{k} f(w_x) \|f(\mathbf{w})\|_1 \frac{p_x}{q_x}. \end{aligned}$$

The upper bound on the variance for a domain  $H$  is:

$$\frac{1}{k} \left( \sum_{x \in H} f(w_x) \right) \|f(\mathbf{w})\|_1 \max_{x \in H} \frac{p_x}{q_x}. \quad (4)$$

Thus for any  $H$ , the variance bound (4) is larger by the benchmark bound (3) by at most a factor of  $h(\mathbf{p}, \mathbf{q})$ .  $\square$

Note that the inaccuracy in the probabilities (using  $\mathbf{q}$  instead of  $\mathbf{p}$ ) is compensated for by a larger sample size without compromising accuracy of the estimate.

*Remark 2.2.* Overhead bounds accumulate multiplicatively: If sampling according to  $\mathbf{q}$  emulates a sample by  $\mathbf{p}$  and a sample by  $\mathbf{q}'$  emulates a sample by  $\mathbf{q}$ , then a sample by  $\mathbf{q}'$  emulates a sample by  $\mathbf{p}$  with overhead  $h(\mathbf{q}', \mathbf{p}) \leq h(\mathbf{q}', \mathbf{q})h(\mathbf{q}, \mathbf{p})$ .

*Remark 2.3.* The emulation overhead can be interpreted as providing an upper bound over all possible estimation tasks that the emulated sample could be used for. This definition of overhead is tight if we wish to transfer guarantees for all subsets  $H$ : Consider a (subset that is a single) key  $x$  and sample size  $k$  such that  $q_x \leq p_x \ll 1/k$ . The variance when sampling according to  $\mathbf{q}$  is  $\approx (f(w_x))^2(1/(kq_x) - 1) \approx (f(w_x))^2/(kq_x) = \frac{1}{k} f(w_x) \frac{f(w_x)}{p_x} \frac{p_x}{q_x} = \frac{1}{k} f(w_x) \|f(\mathbf{w})\|_1 \frac{p_x}{q_x}$ . This is a factor of  $p_x/q_x$  larger than the variance when sampling according to  $p_x = f(w_x)/\|f(\mathbf{w})\|_1$ , which is  $\approx (f(w_x))^2/(kp_x) = \frac{1}{k} f(w_x) \|f(\mathbf{w})\|_1$ .

**Overhead for estimating  $f$ -statistics.** If we are only interested in estimates of the full  $f$ -statistics  $\|f(\mathbf{w})\|_1$ , the overhead reduces to the expected ratio  $\mathbb{E}_{x \sim \mathbf{p}}[p_x/q_x]$  instead of the maximum ratio.

**Corollary 2.4.** *Let  $\mathbf{p}$  be the base pps probabilities for  $f(\mathbf{w})$ . Consider weighted sampling of size  $k$  according to  $\mathbf{q}$ . Then,*

$$\text{Var}[\|\widehat{f(\mathbf{w})}\|_1] \leq \frac{\|f(\mathbf{w})\|_1^2}{k} \sum_x p_x \cdot \frac{p_x}{q_x} = \frac{\|f(\mathbf{w})\|_1^2}{k} \mathbb{E}_{x \sim \mathbf{p}} \left[ \frac{p_x}{q_x} \right].$$

**Multi-objective emulation.** For  $h \geq 1$  and a function of frequency  $f$ , there is a family  $F$  of functions so that a weighted sample according to  $f$  emulates a weighted sample for every  $g \in F$  with overhead  $h$ . A helpful closure property of such  $F$  is the following:

**Lemma 2.5.** *(Cohen, 2015)  $F$  is closed under nonnegative linear combinations. That is, if  $\{f_i\} \subset F$  and  $a_i \geq 0$ , then  $\sum_i a_i f_i \in F$ .*

## 2.1. Off-the-Shelf Composable Sampling Sketches

We consider off-the-shelf use of known polylogarithmic-size sampling sketches. These sketches are designed to provide statistical guarantees on accuracy (bounds on the variance) with respect to specific functions of frequencies and all frequency distributions. The samples still provide unbiased estimates of statistics with respect to any function of frequency. We study the estimates provided by these off-the-shelf sketches through the lens of combinations: instead of considering a particular function of frequency and all frequency distribution, we study the overhead of more general function-frequency combinations.

- (i)  $\ell_1$  sampling without replacement. A ppswor sketch (Cohen et al., 2012) (that builds on the aggregated scheme (Rosén, 1997; Cohen & Kaplan, 2008)

and related schemes (Gibbons & Matias, 1998; Estan & Varghese, 2002)) of size  $k$  (storing  $k$  keys or hashes of keys) performs perfect without replacement sampling of  $k$  keys according to the weights  $w_x$ . The exact frequencies of sampled keys and corresponding inclusion probabilities can be obtained by a second pass over the dataset. Alternatively, in a single streaming pass we can collect partial counts that can be used with appropriate tailored estimators (Cohen et al., 2012).

- (ii)  $\ell_2$  (and generally  $\ell_p$  samples for  $p \in [0, 2]$ ) with replacement. There are multiple designs based on linear projections (Indyk, 2001; Andoni et al., 2011; McGregor et al., 2016). The currently best asymptotic space bound is  $O(\log \delta^{-1} \log^2 n)$  for a single sampled key, where  $\delta$  is the probability of success of producing the sample (Jayaram & Woodruff, 2018). A without-replacement sample of size  $k$  can be obtained with a sketch with  $O(k \log^2 n \log \delta^{-1})$  bits (assuming the keys are integers between 1 and  $n$ ).<sup>2</sup> We note that on skewed distributions without-replacement sampling is significantly more effective for a fixed sample size. Our understanding is that there is no prior work for sketches of size  $\tilde{O}(k)$  for without-replacement schemes for  $p > 1$  but there is a planned publication for such sketches for  $p \leq 2$  (Cohen et al., 2020).

- (iii) Sampling sketches for capping functions (Cohen, 2018) and more generally concave sublinear functions (Cohen & Geri, 2019). We will also consider a multi-objective sample that emulates all concave sublinear functions of frequency with space overhead  $O(\log n)$  (Cohen, 2018).

In our analysis and experiments we compute the overhead of using the above sketches with respect to certain frequencies and functions. Recall that for each target function of frequency, the overhead is computed with respect to the applicable base pps probabilities  $p_x = \frac{f(w_x)}{\|f(\mathbf{w})\|_1}$ . Consider a frequency vector  $\mathbf{w}$  in nonincreasing order ( $w_i \geq w_{i+1}$ ). The base pps sampling for  $\ell_p$  sampling are simply  $w_i^p / \|\mathbf{w}\|_p^p$ . The base pps probabilities for multi-objective concave sublinear sampling are  $q_i = q'_i / \|\mathbf{q}'\|_1$  where  $q'_i := \frac{w_i}{i w_i + \sum_{j=i+1}^n w_j}$ . These samples emulate sampling for all concave-sublinear functions with overhead  $\|\mathbf{q}'\|_1$ .

## 3. The Advice Model

In this section, we assume that in addition to the input, we are provided with an oracle access to an ‘‘advice’’ model. When presented with a key  $x$ , the advice model returns a

<sup>2</sup>The dependence on  $\delta$  improves with  $k$  but we omit this for brevity.

prediction  $a_x$  for the total frequency of  $x$  in the data. For simplicity, we assume that predictions are the same for all queries with the same key. This model is similar to the model used in a recent paper about frequency estimation (Hsu et al., 2019).

A detailed description of our sketch structure and the corresponding estimators (including proofs) is provided in Appendix B. At a high level, our sampling sketch takes size parameters  $(k_h, k_p, k_u)$ , maintains a set of at most  $k_h + k_p + k_u$  keys, and collects the exact frequencies  $w_x$  for these stored keys. The primary component of the sketch is a weighted-sample-by-advice of size  $k_p$ . Our sketch stores keys according to two additional criteria in order to provide robustness to the prediction quality of the advice:

- The top- $k_h$  keys by advice. This provides tolerance to inaccuracies in the advice for these heaviest keys. Since these keys are included with probability 1, they will not contribute to the error.
- A uniform sample of  $k_u$  keys. This allows keys that are "below average" in their contribution to  $\|f(\mathbf{w})\|_1$  to be represented appropriately in the sample, regardless of the accuracy of the advice. This provides robustness to the accuracy of the advice on these very infrequent keys and ensures they are not undersampled. Moreover, this ensures that all active keys ( $w_x > 0$ ), including those with potentially no advice ( $a_x = 0$ ), have a positive probability of being sampled. This is necessary for unbiased estimation.

We provide an unbiased estimator that smoothly combines the different sketch components and provides the following guarantees:

**Lemma 3.1.** *Suppose the advice model is such that for some  $c_p, c_u \geq 0$  and  $h \geq 0$ , all keys  $x$  that are active ( $w_x > 0$ ) and not in the  $h$  largest advice values of active keys ( $a_x < \{a_y \mid w_y > 0\}_{(n-h+1)}$ ) satisfy*

$$\frac{f(w_x)}{\|f(\mathbf{w})\|_1} \leq \max\left\{c_p \frac{f(a_x)}{\|f(\mathbf{a})\|_1}, c_u \frac{1}{n}\right\}.$$

Then for all  $k \geq 1$ , a sample with  $(k_h, k_p, k_u) = (h, \lceil kc_p \rceil + 2, \lceil kc_u \rceil + 2)$  satisfies the variance bound (3) for all  $H$ .

In particular, if our advice is approximately accurate, say  $f(w_x) \leq f(a_x) \leq c_p \cdot f(w_x)$ , the overhead when sampling by advice is  $c_p$ .

**Corollary 3.2.** *Let  $f$  be such that  $f(w_x) \leq f(a_x) \leq c_p f(w_x)$  then for all  $k \geq 1$ , with sample size  $(k_h, k_p, k_u) = (0, \lceil kc_p \rceil + 2, 0)$  we have  $\text{Var} \left[ \|f(\widehat{\mathbf{w}})\|_1 \right] \leq \frac{1}{k} \|f(\mathbf{w})\|_1^2$ .*

**Experiments.** We evaluate the effectiveness of "sampling by advice" for estimating the frequency moments with

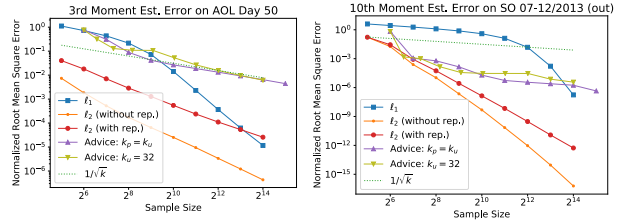


Figure 1. Estimating the third moment on the AOL dataset (with learned advice) and the tenth moment on the Stack Overflow dataset (with past frequencies).

$p = 3, 7, 10$  on datasets from (Pass et al., 2006; Paranjape et al., 2017) (see details in Appendix A). We use advice models from prior work (Hsu et al., 2019) based on a machine learning algorithms applied to past data and advice based directly on frequencies in past data. Some representative results are reported in Figure 1 and additional results and more details are provided in Appendix C. The results reported for sampling by advice are with  $k_h = 0$  and two choices of balance between the ppswor sample based on the advice and the uniform sample:  $k_p = k_u$  and  $k_u = 32$ . We also report performance of ppswor ( $\ell_1$  sampling without replacement),  $\ell_2$  sampling (with and without replacement), and the benchmark upper bound ( $1/\sqrt{k}$ ).

We observe that for the third moment, sampling by advice did not perform significantly better than ppswor (and sometimes performed even worse). For higher moments, however, sampling by advice performed better than ppswor when the sample size is small. With replacement  $\ell_2$  sampling was more accurate than advice (without replacement  $\ell_2$  sampling performs best). Our analysis in the next section explains the perhaps surprisingly good performance of  $\ell_1$  and  $\ell_2$  sampling schemes.

## 4. Frequencies-Functions Combinations

In this section, we analyze performance for inputs that comes from restricted families of frequency distributions  $W$ . Restricting the family  $W$  of possible inputs allows us to extend the family  $F$  of functions of frequency that can be efficiently emulated with a small overhead.

Specifically, we will consider sampling sketches and corresponding combinations  $(W, F, h)$  of frequency vectors  $W$ , functions of frequency  $F$ , and overhead  $h \geq 1$  so that for every frequency distribution  $w \in W$  and frequency function  $f \in F$ , our sampling sketch emulates a weighted sample of  $f(\mathbf{w})$  with overhead  $h$ . We will say that our sketch supports the combination  $(W, F, h)$ .

Recall that emulation with overhead  $h$  means that a sampling sketch of size  $h\epsilon^{-2}$  (i.e., holding this number of keys or hashes of keys) provides estimates with NRMSE  $\epsilon$  for  $f(\mathbf{w})$

for any  $f \in F$  and that these estimates are concentrated in the Chernoff bound sense. Moreover, for any  $f \in F$  we can estimate statistics of the form (1) with the same guarantees on accuracy as provided by a dedicated weighted sample according to  $f$ .

We study combinations supported by off-the-shelf sampling schemes that can be implemented with small (e.g. polylogarithmic) size sketches as detailed in Section 2.1. We will use the notation  $\mathbf{w} = \{w_i\}$  for dataset frequencies, where  $w_i$  is the frequency of the  $i$ th most frequent key.

We reports results of experiments on datasets listed in Table 1 with details provided in in Appendix A.

#### 4.1. Emulating an $\ell_p$ Sample by an $\ell_q$ Sample

We express the overhead of emulating  $\ell_p$  sampling by  $\ell_q$  sampling ( $p \geq q$ ) in terms of the properties of the frequency distribution. Recall that  $\ell_p$  sampling (and estimates of  $p$ -th frequency moments) can be implemented with polylogarithmic size sketches for  $p \leq 2$  but requires polynomial size sketches in the worst case when  $p > 2$ .

**Lemma 4.1.** *Consider a dataset with frequencies  $\mathbf{w}$  (in nonincreasing order). For  $p \geq q$ , the overhead of emulating  $\ell_p$  sampling by  $\ell_q$  sampling is bounded by*

$$\frac{\left\| \frac{\mathbf{w}}{w_1} \right\|_q^q}{\left\| \frac{\mathbf{w}}{w_1} \right\|_p^p}. \quad (5)$$

*Proof.* The sampling probabilities for key  $i$  under  $\ell_p$  sampling and  $\ell_q$  sampling are  $\frac{w_i^p}{\|\mathbf{w}\|_p^p}$  and  $\frac{w_i^q}{\|\mathbf{w}\|_q^q}$ , respectively. Then, the overhead of emulating  $\ell_p$  sampling by  $\ell_q$  sampling is

$$\max_i \frac{w_i^p / \|\mathbf{w}\|_p^p}{w_i^q / \|\mathbf{w}\|_q^q} = \max_i w_i^{p-q} \cdot \frac{\|\mathbf{w}\|_q^q}{\|\mathbf{w}\|_p^p} = w_1^{p-q} \cdot \frac{\|\mathbf{w}\|_q^q}{\|\mathbf{w}\|_p^p} = \frac{\left\| \frac{\mathbf{w}}{w_1} \right\|_q^q}{\left\| \frac{\mathbf{w}}{w_1} \right\|_p^p}.$$

□

We can obtain a (weaker) upper bound on the overhead, expressed only in terms of  $q$ , that applies to all  $p \geq q$ :

**Corollary 4.2.** *The overhead of emulating  $\ell_p$  sampling using  $\ell_q$  sampling (for any  $p \geq q$ ) is at most  $\left\| \frac{\mathbf{w}}{w_1} \right\|_q^q$ .*

*Proof.* For any set of frequencies  $\mathbf{w}$ , the normalized norm  $\left\| \frac{\mathbf{w}}{w_1} \right\|_p^p$  is non-increasing with  $p$  and is at least 1. Therefore, the overhead (5) is

$$\left\| \frac{\mathbf{w}}{w_1} \right\|_q^q \Big/ \left\| \frac{\mathbf{w}}{w_1} \right\|_p^p \leq \left\| \frac{\mathbf{w}}{w_1} \right\|_q^q.$$

□

**Remark 4.3.** Emulation works when  $p \geq q$ . When  $q > p$ , the maximum in the overhead bound (see proof of Lemma 4.1) is incurred on the least frequent key, with frequency  $w_n$ . We therefore get a bound of  $\left\| \frac{\mathbf{w}}{w_n} \right\|_q^q / \left\| \frac{\mathbf{w}}{w_n} \right\|_p^p$  and Corollary 4.2 does not apply.

#### 4.2. Frequency Distributions with a Heavy Hitter

We show that for distributions with an  $\ell_q$  heavy hitter,  $\ell_q$  sampling emulates  $\ell_p$  sampling for all  $p \geq q$  with a small overhead.

**Definition 4.4.** *Consider frequencies  $\mathbf{w}$ . An  $\ell_q$   $\phi$ -heavy hitter is defined to be a key such that  $w_i^q \geq \phi \cdot \|\mathbf{w}\|_q^q$ .*

We rephrase Corollary 4.2 in terms of a presence of a heavy hitter:

**Corollary 4.5.** *Let  $\mathbf{w}$  be a frequency vector with a  $\phi$ -heavy hitter under  $\ell_q$ . Then for  $p \geq q$ , the overhead of using  $\ell_q$  sample to emulate an  $\ell_p$  sample is at most  $1/\phi$ .*

*Proof.* If there is an  $\ell_q$  heavy hitter then then the most frequent key (the key with frequency  $w_1$ ) must be a heavy hitter. From the definition of a heavy hitter,  $\left\| \frac{\mathbf{w}}{w_1} \right\|_q^q \leq \frac{1}{\phi}$ , and we get the desired bound on the overhead. □

We are now ready to specify combinations  $(W, F, h)$  of frequency vectors  $W$ , functions of frequency  $F$ , and overhead  $h \geq 1$  that are supported by  $\ell_q$  sampling.

**Theorem 4.6.** *For any  $q > 0$  and  $\phi \in (0, 1]$ , an  $\ell_q$ -sample supports the combination*

$$\begin{aligned} W &:= \{\mathbf{w} \text{ with an } \ell_q \text{ } \phi\text{-heavy hitter}\} \\ F &:= \overline{\{f(\mathbf{w}) = w^p \mid p \geq q\}}_+ \\ h &:= 1/\phi, \end{aligned}$$

where the notation  $\overline{F}_+$  is the closure of a set  $F$  of functions under nonnegative linear combinations.

*Proof.* The claim for functions  $f(\mathbf{w}) = w^p$  is immediate from Corollary 4.5. The claim for the nonnegative convex closure of these function is a consequence of Lemma 2.5. □

In particular, if the input distribution has an  $\ell_q$   $\phi$ -heavy hitter then  $\ell_q$  sampling of size  $\varepsilon^{-2}/\phi$  emulates an  $\ell_p$  sampling of size  $\varepsilon^{-2}$  for any  $p > q$ .

Table 1 reports properties and the relative  $\ell_1$  and  $\ell_2$  weights of the most frequent key for our datasets. We can see that the most frequent key is a heavy hitter with  $1/\phi \leq 21$  for  $\ell_2$  and  $1/\phi \leq 625$  for  $\ell_1$  which gives us upper bounds on the overhead of emulating any  $\ell_p$  sample ( $p \geq 2$ ). Table 3

Table 1. Datasets

Dataset	$n/10^6$	$\ln n$	$\ell_1$ HH	$\ell_2$ HH	$\alpha$
SO in	2.23	14.6	0.0016, 0.0010	0.053, 0.022	1.48
SO out	2.30	14.6	0.0015, 0.0009	0.107, 0.036	1.38
AOL	10.15	16.1	0.0275, 0.0091	0.827, 0.091	0.77
CAIDA	1.07	13.9	0.0033, 0.0032	0.048, 0.046	1.35
UGR	79.38	18.2	0.1118, 0.0401	0.850, 0.109	1.35

reports (for  $p = 3, 7, 10$ ) the overhead of emulating the respective  $\ell_p$  sample and the (smaller) overhead of estimating the  $p$ th moment. We can see that high moments can be estimated well from  $\ell_2$  and with a larger overhead from  $\ell_1$  samples.

**Certified emulation.** The quality guarantees of a combination  $(W, F, h)$  are provided when  $w \in W$ . In practice, however, we may compute samples of arbitrary dataset frequencies  $w$ . Conveniently, we are able to test the validity of emulation by considering the most frequent key in the sample: For an  $\ell_q$  sample of size  $k$  we can compute  $r \leftarrow \max_{x \in S} w_x^q / \|w\|_q^q$  and certify that our sample emulates  $\ell_p$  samples ( $p > q$ ) of size  $kr$ . If  $kr$  is small, then we do not provide meaningful accuracy but otherwise we can certify the emulation with sample size  $kr$ . When the input  $w$  has an  $\ell_q$   $\phi$ -heavy hitter then an  $\ell_q$  sample of size  $k$  will include it with probability at least  $1 - e^{-k\phi}$  and the result will be certified. Note that the result can only be certified if there is a heavy hitter.

**Tradeoff between  $W$  and  $F$ .** If  $w$  has an  $\ell_q$   $\phi$ -heavy hitter then it has an  $\ell_p$   $\phi$ -heavy hitter for every  $p \geq q$ . This means that for moments  $p \geq 2$ , an  $\ell_2$  sample supports a larger set  $W$  of frequencies than an  $\ell_1$  sample, including also those with an  $\ell_2$   $\phi$ -heavy hitter but not an  $\ell_1$   $\phi$ -heavy hitter. The  $\ell_1$  sample however supports a larger family  $F$  that includes moments with  $p \in [1, 2)$ . Note that for a fixed overhead and  $\ell_q$  sampling, the set  $F$  of supported functions decreases with  $q$  whereas  $W$  increases with  $q$ .

### 4.3. Zipfian and sub-Zipfian Frequencies

Zipf distributions are a very common model for frequency distributions in practice. We explore supported combinations with frequencies that are (approximately) Zipf.

**Definition 4.7.** We say that the frequencies  $w$  where  $\|w\|_0 = n$  are  $\text{Zipf}[\alpha, n]$  (Zipf with parameter  $\alpha$ ) if for all  $i$ ,  $w_i/w_1 = i^{-\alpha}$ .

Values  $\alpha \in [1, 2]$  are common in practice. The best-fit Zipf parameter for the datasets we studied is reported in Table 1 and the frequency distribution (sorted by rank) is shown in Figure 2. We can see that our datasets are approximately Zipf (which would be an approximate straight line) and for all but one we have  $\alpha \in [1.3, 1.5]$ .

Table 2. Supported combinations for  $\text{subZipf}[\alpha, c, n]$  frequencies

method	$W$	$F$	Overhead
$\ell_1$ sampling	$\{\alpha \geq 2\}$	$\{f(w) = w^p \mid p \geq 1\}_+$	$1.65c$
$\ell_1$ sampling	$\{\alpha \geq 1\}$	$\{f(w) = w^p \mid p \geq 1\}_+$	$(1 + \ln n)c$
$\ell_2$ sampling	$\{\alpha \geq 1\}$	$\{f(w) = w^p \mid p \geq 2\}_+$	$1.65c^2$
$\ell_2$ sampling	$\{\alpha \geq 1/2\}$	$\{f(w) = w^p \mid p \geq 2\}_+$	$(1 + \ln n)c^2$

We now define a broader class of approximately Zipfian distributions.

**Definition 4.8.** Frequencies  $w$  are  $\text{subZipf}[\alpha, c, n]$  if for all  $i$ ,  $\frac{w_i}{w_1} \leq ci^{-\alpha}$ .<sup>3</sup>

Note that  $\text{Zipf}[\alpha, n]$  is sub-Zipfian with the same  $\alpha$  and  $c = 1$ . We show that sub-Zipf frequencies (and in particular Zipf frequencies) have heavy hitters:

**Lemma 4.9.** For  $\text{subZipf}[\alpha, c, n]$  frequencies, and  $q$  such that  $q\alpha \geq 1$ , the frequency vector has an  $\ell_q$   $c^q \frac{1}{H_{n,\alpha}}$ -heavy hitter, where  $H_{n,\alpha} := \sum_{i=1}^n i^{-\alpha}$  are the generalized harmonic numbers.

*Proof.* We use Corollary 4.2 to express the overhead (5) for  $\text{subZipf}[\alpha, c, n]$  frequencies

$$\left\| \frac{w}{w_1} \right\|_q^q \leq c^q \sum_{i=1}^n i^{-\alpha q} = c^q H_{n,q\alpha} . \quad (6)$$

□

Table 2 lists supported combinations that include these approximately Zipfian distributions.

**Lemma 4.10.** The combinations shown in Table 2 are supported by  $\ell_1$  and  $\ell_2$  samples.

*Proof.* We use Lemma 4.9 and Theorem 4.6. Recall that when  $\alpha = 1$ , the harmonic sum is  $H_{n,1} \leq 1 + \ln n$ . For  $\alpha > 1$ ,  $H_{n,\alpha} \leq \zeta(\alpha)$ , where the Zeta function  $\zeta(\alpha) := \sum_{i=1}^{\infty} i^{-\alpha}$ . The Zeta function is decreasing with  $\alpha$ , defined for  $\alpha > 1$  with an asymptote at  $\alpha = 1$ , and is at most 1.65 for  $\alpha \geq 2$ .

When  $q\alpha \geq 2$ , the overhead is  $\leq c^q \zeta(q\alpha) \leq 1.65c^q$ . When  $q\alpha = 1$  the overhead is at most  $(1 + \ln n)c^q$  and when  $q\alpha > 1$  we can bound it by  $\min\{1 + \ln n, \zeta(q\alpha)\}c^q$ . □

We see that for these approximately Zipf distributions,  $\ell_1$  or  $\ell_2$  samples emulate  $\ell_p$  samples with small overhead.

### 4.4. Experiments on Estimate Quality

The overhead factors reported in Table 3 are in a sense worst-case upper bounds (for the dataset frequencies). Figure 4

<sup>3</sup>This is a slight abuse of notation since the parameters do not fully specify a distribution



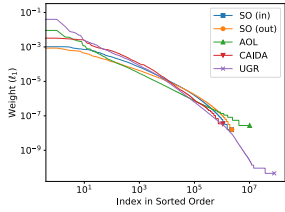


Figure 2. Freq. by rank

reports the actual estimation error (normalized root mean square error) for high moments for representative datasets as a function of sample size. The estimates are with ppswor ( $\ell_1$  sample with replacement) and  $\ell_2$  samples with and without replacement. Additional results are reported in Appendix E. We observe that actual accuracy is significantly better than even the benchmark bounds.

Finally we consider estimating the full distribution of frequencies, that is, the curve that relates frequency of keys to their rank. We do this by estimating the actual rank of each key in the sample (using an appropriate threshold function of frequency). Representative results are reported in Figure 3 for ppswor and for with-replacement  $\ell_2$  sampling (additional results are reported in Appendix E). We used a sample of size  $k = 32$  or  $k = 1024$  for each set of estimates. We observe that generally the estimates are fairly accurate even with a small sample size (despite threshold function requiring large sketches in the worst case). We see that  $\ell_2$  samples are accurate for the frequent keys but often have no representatives from the tail whereas the without replacement  $\ell_1$  samples are more accurate on the tail.

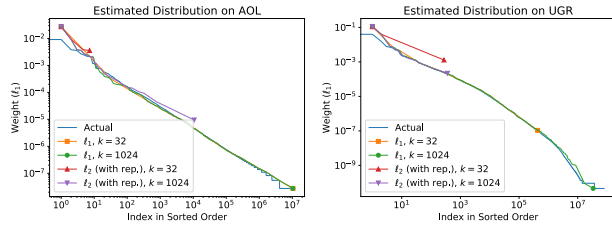


Figure 3. Actual and estimated distribution of frequency by rank. Estimates for ppswor and with-replacement  $\ell_2$  sampling and sample sizes  $k = 32, 1024$ .

#### 4.5. Worst-Case Bound on Overhead

The overhead (5) of  $\|\mathbf{w}/w_1\|_2^2/\|\mathbf{w}/w_1\|_p^2$  is the space factor increase needed for an  $\ell_2$  sample to emulate an  $\ell_p$  sample on the frequencies  $\mathbf{w}$  (and estimate accurately  $p$ th moments). A natural question is whether there is a better way to emulate an  $\ell_p$  sampling with a polylogarithmic size sampling sketch. The following shows that in a sense an  $\ell_2$  sample is the best we can do:

Table 3. Overhead

Dataset	$\ell_1$ max overhead (expected overhead)			$\ell_2$ max overhead (expected overhead)			concave universal
	3rd	10th	universal	3rd	10th	universal	
SO in	124.30 (42.76)	600.44 (577.57)	624.58	3.74 (1.90)	18.04 (17.36)	$1.25 \times 10^5$	1672.60
SO out	299.80 (155.32)	677.56 (673.45)	681.72	4.12 (2.58)	9.30 (9.25)	$4.20 \times 10^4$	1628.36
AOL	34.81 (33.45)	36.38 (36.37)	92.92	1.16 (1.12)	1.21 (1.21)	$2.94 \times 10^5$	170.84
CAIDA	31.23 (18.73)	90.66 (56.20)	301.28	2.16 (1.57)	6.28 (4.03)	$2.65 \times 10^5$	846.15
UGR	8.52 (8.17)	8.95 (8.95)	772.83	1.12 (1.09)	1.18 (1.18)	$1.89 \times 10^{11}$	143.94

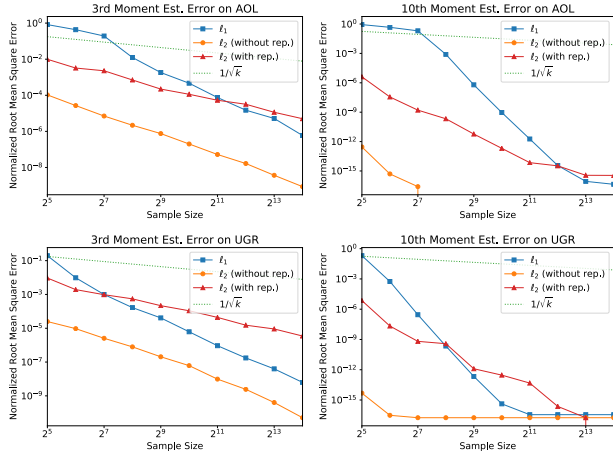


Figure 4. Estimating 3rd and 10th moments on various datasets from ppswor ( $\ell_1$  without replacement) and  $\ell_2$  samples (with and without replacement). The error is averaged over 50 repetition.

**Lemma 4.11.** An  $\ell_2$  sample with overhead  $n^{1-2/p}$  can emulate an  $\ell_p$  sample for any frequency vector  $\mathbf{w}$  with support  $\|\mathbf{w}\|_0 = n$ .

*Proof.*  $\max_{\mathbf{w}} \frac{\|\mathbf{w}/w_1\|_2^2}{\|\mathbf{w}/w_1\|_p^2} \leq n^{1-2/p}$ .  $\square$

This matches upper bounds on sketch size attained with dedicated sketches for  $p$ th-moment estimation (Indyk & Woodruff, 2005; Andoni et al., 2011) and the worst-case lower bound of  $\tilde{\Omega}(n^{1-2/p})$  (Alon et al., 1999; Li & Woodruff, 2013). Interestingly, the worst case distributions that establish that bound are those where the most frequent key is  $\ell_p$  heavy but not  $\ell_2$  heavy.

## 5. Universal Samples

We study combinations where the emulation is *universal*, that is,  $F$  includes the set  $M$  of all monotone non-decreasing functions of frequency. Interestingly, there are sampling probabilities that provide universal emulation for any  $\mathbf{w}$ :

**Lemma 5.1.** (Cohen, 2015) Consider the probabilities  $\mathbf{q}$  where the  $i$ th most frequent key has  $q_i = \frac{1}{iH_n}$ . Then a

weighted sample by  $\mathbf{q}$  is a universal emulator with overhead at most  $H_n$ .

*Proof.* Consider a monotone non-decreasing  $f$  with respective pps probabilities  $\mathbf{p}$ . By definition, for the  $i$ th most frequent key,  $p_i = \frac{f(w_i)}{\|f(\mathbf{w})\|_1} \leq \frac{1}{i}$ . Therefore,  $p_i/q_i \leq 1/i(iH_n) = H_n q_i$ .  $\square$

This universal sampling, however, can not be implemented with small (polylogarithmic size) sketches. This because  $M$  includes functions that require large (polynomial size) sketches such as thresholds ( $I_{w \geq T}$  for some  $T$ ) and high moments ( $p > 2$ ). We therefore aim for small sketches that provide universal emulation to a restricted  $W$ .

For particular sampling probabilities  $\mathbf{q}$  and frequencies  $\mathbf{w}$  we consider the *universal emulation overhead* to be the overhead factor that will allow the sample to emulate weighted sampling with respect to  $f(\mathbf{w})$  for any  $f \in M$ .

$$\max_{f \in M} \max_i f(w_i) / (\|f(\mathbf{w})\|_1 q_i) \quad (7)$$

Interestingly, the universal emulation overhead of  $\mathbf{q}$  does not depend on the particular  $\mathbf{w}$ .

**Lemma 5.2.** *The universal emulation overhead of  $\mathbf{q}$  is*

$$\max_i 1/(iq_i)$$

and is always at least  $H_n$ . This is tight even when  $W$  contains a single  $\mathbf{w}$ , as long as frequencies are distinct ( $w_i > w_{i+1}$  for all  $i$ ).

*Proof.* Consider  $\mathbf{w}$ . We have

$$\begin{aligned} & \max_{f \in M} \max_i f(w_i) / (\|f(\mathbf{w})\|_1 q_i) \\ &= \max_i \max_{f \in M} f(w_i) / (\|f(\mathbf{w})\|_1 q_i) \leq \max_i 1/(iq_i). \end{aligned}$$

The last inequality follows because  $f \in M$ , for all  $j < i$  we must have  $f(w_j) \geq f(w_i)$ . Therefore  $f(w_i) / \|f(\mathbf{w})\|_1 \leq 1/i$ . This is maximized for the threshold function at  $w_i$  and equality holds (if  $w_{i+1} < w_i$ ).  $\square$

We can similarly consider for sampling probabilities  $\mathbf{q}$  the *universal estimation overhead* which is the overhead needed for estimating all (full) monotone  $f$ -statistics. As discussed in Section 2, the estimation is a weaker requirement than emulation (only applies to the full statistics) and hence for any particular  $\mathbf{q}$  the estimation overhead can be lower than the emulation overhead. The estimation overhead, however, is still at least  $H_n$ .

**Lemma 5.3.** *The universal estimation overhead for estimating all monotone  $f$ -statistics for  $\mathbf{q}$  is*

$$\max_i \frac{1}{i^2} \sum_{j=1}^i \frac{1}{q_j}.$$

*Proof.* The overhead with frequencies  $\mathbf{w}$  is

$$\max_{f \in M} \sum_i \frac{f(w_i)^2}{\|f(\mathbf{w})\|_1^2} \frac{1}{q_i}. \quad (8)$$

It suffices to consider  $f$  that are threshold functions. The expression for the threshold at  $w_i$  has  $f(w_j) / \|f(\mathbf{w})\|_1 = 1/i$  for  $j \leq i$  and 0 otherwise. We get that the sum is  $\frac{1}{i^2} \sum_{j=1}^i \frac{1}{q_j}$ . The claim follows from taking the maximum over all threshold functions.  $\square$

In our context, the probabilities  $\mathbf{q}$  are not something we directly control but rather emerge as an artifact of applying a certain sampling scheme to a dataset with certain frequencies  $\mathbf{w}$ . We will explore the universal overhead of  $\mathbf{q}$  we obtain when applying off-the-shelf schemes (see Section 2.1) to Zipf frequencies and to our datasets.

For Zipf $[\alpha]$  frequencies ( $\alpha \geq 1/2$ ),  $\ell_p$  sampling with  $p = 1/\alpha$  is a universal emulator with (optimal) overhead  $H_n$ . Interestingly, for  $\alpha \geq 1/2$ , this is attained by  $\ell_p$  sampling with  $p \leq 2$ , which has polylogarithmic size sketches. Note that we match here a different  $\ell_p$  sample for each possible Zipf parameter  $\alpha$  of the data frequencies. A sampling scheme that emulates  $\ell_p$  sampling for a range  $[p_1, p_2]$  of  $p$  values with some overhead  $h$  will be a universal emulator with overhead  $hH_n$  for Zipf $[\alpha]$  for  $\alpha \in [1/p_2, 1/p_1]$  (see Remark 2.2). One such sampling scheme with polylogarithmic-size sketches was provided in (Cohen, 2018; 2015). The sample emulates all concave sublinear functions that include capping functions  $f(w) = \min\{w, T\}$  for  $T > 0$  and low moments with  $p \in [0, 1]$  with  $O(\log n)$  overhead.

Table 3 reports the universal overhead on our datasets with  $\ell_1, \ell_2$ , and multi-objective concave-sublinear sampling probabilities. We observe that while  $\ell_2$  sampling emulates high moments extremely well, its universal overhead is very large due to poor emulation of “slow growth” functions. The better universal overhead of  $\ell_1$  and concave-sublinear samples is  $h \in [143, 700]$  and is practically meaningful as it is in the regime where  $h\varepsilon^{-2} \ll n$ .

The universal overhead was computed using Lemma 5.2 with respect to base pps probabilities for the sampling schemes as described in Section 2.1.

We next express a condition on the frequency distribution under which a multi-objective concave-sublinear sample provides a universal emulator. The condition is that for all  $i$ , the weight of the  $i$ th most frequent key is at least  $c/i$  times the weight of the tail from  $i$ .

**Lemma 5.4.** *Let  $\mathbf{w}$  be such that  $\min_i \frac{iw_i}{\sum_{j=i+1}^n w_j} \geq c$ . Then a sample that emulates all concave-sublinear functions with overhead  $h'$  is a universal emulator for  $\mathbf{w}$  with overhead  $h'(1 + 1/c)H_n$ .*

*Proof.* Consider the  $i$ th largest frequency  $w_i$ . The capping function  $f(w) := \min\{w, w_i\}$  maximizes the sampling probability of key  $i$ , which is

$$\frac{f(w_i)}{\sum_j f(w_j)} = \frac{w_i}{iw_i + \sum_{j>i} w_j} \leq \frac{w_i}{iw_i + \frac{i}{c}w_i} = \frac{1}{i(1+1/c)}.$$

□

Interestingly, for high moments to be “easy” it sufficed to have a heavy hitter. For universal emulation we need to bound from below the relative weight of each key with respect to the remaining tail.

## 6. Conclusion

We propose a framework where performance and statistical guarantees of sampling sketches are analyzed in terms of supported *frequencies-functions combinations*. We demonstrate analytically and empirically that sketches originally designed to sample according to “easy” functions of frequency on “hard” frequency distributions turn out to be accurate for sampling according to “hard” functions of frequency on “practical” frequency distributions. In particular, on “practical” distributions we can accurately approximate high frequency moments ( $p > 2$ ) and the rank versus frequency distribution using small composable sketches.

## Acknowledgments

We are grateful to the authors of (Hsu et al., 2019), especially Chen-Yu Hsu and Ali Vakilian, for sharing their data, code, and predictions with us. We thank Ravi Kumar and Robert Krauthgamer for helpful discussions. The work of Rasmus Pagh was supported by Investigator Grant 16582, Basic Algorithms Research Copenhagen (BARC), from the VILLUM Foundation.

## References

- Aamand, A., Indyk, P., and Vakilian, A. (Learned) frequency estimation algorithms under zipfian distribution. *CoRR*, abs/1908.05198, 2019. URL <http://arxiv.org/abs/1908.05198>.
- Agarwal, P. K., Cormode, G., Huang, Z., Phillips, J. M., Wei, Z., and Yi, K. Mergeable summaries. *ACM Trans. Database Syst.*, 38(4):26:1–26:28, 2013. doi: 10.1145/2500128. URL <https://doi.org/10.1145/2500128>.
- Alon, N., Matias, Y., and Szegedy, M. The space complexity of approximating the frequency moments. *Journal of Computer and System Sciences*, 58(1):137 – 147, 1999. ISSN 0022-0000. doi: <https://doi.org/10.1006/jcss.1997.1545>.
- Alon, N., Gibbons, P. B., Matias, Y., and Szegedy, M. Tracking join and self-join sizes in limited storage. *Journal of Computer and System Sciences*, 64(3):719–747, 2002.
- Andoni, A., Krauthgamer, R., and Onak, K. Streaming algorithms via precision sampling. In *IEEE 52nd Annual Symposium on Foundations of Computer Science, FOCS 2011, Palm Springs, CA, USA, October 22-25, 2011*, 2011. URL <https://doi.org/10.1109/FOCS.2011.82>.
- Braverman, V. and Ostrovsky, R. Zero-one frequency laws. In *STOC*. ACM, 2010.
- CAIDA. The caida ucsd anonymized internet traces 2016 - 2016/01/21 13:29:00 utc. <https://www.caida.org/data/passive/passive.dataset.xml>, 2016.
- Chao, M. T. A general purpose unequal probability sampling plan. *Biometrika*, 69(3):653–656, 1982.
- Charikar, M., Chen, K., and Farach-Colton, M. Finding frequent items in data streams. In *Proceedings of the International Colloquium on Automata, Languages and Programming (ICALP)*, pp. 693–703, 2002.
- Cohen, E. Multi-objective weighted sampling. In *2015 Third IEEE Workshop on Hot Topics in Web Systems and Technologies (HotWeb)*, pp. 13–18, 2015. doi: 10.1109/HotWeb.2015.8. Full version: <https://arxiv.org/abs/1509.07445>.
- Cohen, E. Stream sampling framework and application for frequency cap statistics. *ACM Trans. Algorithms*, 14(4):52:1–52:40, 2018. ISSN 1549-6325. doi: 10.1145/3234338.
- Cohen, E. and Geri, O. Sampling sketches for concave sublinear functions of frequencies. In *NeurIPS*, 2019.
- Cohen, E. and Kaplan, H. Tighter estimation using bottom-k sketches. In *Proceedings of the 34th VLDB Conference*, 2008. URL <http://arxiv.org/abs/0802.3448>.
- Cohen, E., Kaplan, H., and Sen, S. Coordinated weighted sampling for estimating aggregates over multiple weight assignments. *VLDB*, 2, 2009. full version: <http://arxiv.org/abs/0906.4560>.
- Cohen, E., Duffield, N., Lund, C., Thorup, M., and Kaplan, H. Efficient stream sampling for variance-optimal estimation of subset sums. *SIAM J. Comput.*, 40(5), 2011.
- Cohen, E., Cormode, G., and Duffield, N. Don’t let the negatives bring you down: Sampling from streams of signed updates. In *Proc. ACM SIGMETRICS/Performance*, 2012.

- Cohen, E., Pagh, R., and Woodruff, D. P. Personal communication, 2020.
- Cormode, G. and Muthukrishnan, S. An improved data stream summary: The count-min sketch and its applications. *J. Algorithms*, 55(1), 2005.
- Duffield, N., Thorup, M., and Lund, C. Priority sampling for estimating arbitrary subset sums. *Journal of the ACM*, 54(6), 2007.
- Eden, T., Ron, D., and Seshadhri, C. Sublinear time estimation of degree distribution moments: The arboricity connection. *SIAM J. Discrete Math.*, 33, 2019. URL <https://doi.org/10.1137/17M1159014>.
- Estan, C. and Varghese, G. New directions in traffic measurement and accounting. In *SIGCOMM*. ACM, 2002.
- Flajolet, P. and Martin, G. N. Probabilistic counting algorithms for data base applications. *Journal of Computer and System Sciences*, 31:182–209, 1985.
- Flajolet, P., Fusy, E., Gandouet, O., and Meunier, F. Hyperloglog: The analysis of a near-optimal cardinality estimation algorithm. In *Analysis of Algorithms (AofA)*. DMTCS, 2007.
- Frieze, A., Kannan, R., and Vempala, S. Fast monte-carlo algorithms for finding low-rank approximations. *J. ACM*, 51(6), 2004. URL <https://doi.org/10.1145/1039488.1039494>.
- Gibbons, P. and Matias, Y. New sampling-based summary statistics for improving approximate query answers. In *SIGMOD*. ACM, 1998.
- Horvitz, D. G. and Thompson, D. J. A generalization of sampling without replacement from a finite universe. *Journal of the American Statistical Association*, 47(260):663–685, 1952. ISSN 01621459. URL <http://www.jstor.org/stable/2280784>.
- Hsu, C.-Y., Indyk, P., Katabi, D., and Vakilian, A. Learning-based frequency estimation algorithms. In *International Conference on Learning Representations*, 2019. URL <https://openreview.net/forum?id=r1lohoCqY7>.
- Indyk, P. Stable distributions, pseudorandom generators, embeddings and data stream computation. In *Proc. 41st IEEE Annual Symposium on Foundations of Computer Science*, pp. 189–197. IEEE, 2001.
- Indyk, P. and Woodruff, D. P. Optimal approximations of the frequency moments of data streams. In *STOC*, pp. 202–208. ACM, 2005.
- Indyk, P., Vakilian, A., and Yuan, Y. Learning-based low-rank approximations. In *NeurIPS*, 2019.
- Jayaram, R. and Woodruff, D. P. Perfect lp sampling in a data stream. In *FOCS*, 2018. URL <https://doi.org/10.1109/FOCS.2018.00058>.
- Jiang, T., Li, Y., Lin, H., Ruan, Y., and Woodruff, D. P. Learning-augmented data stream algorithms. In *International Conference on Learning Representations*, 2020. URL <https://openreview.net/forum?id=HyxJ1xBYDH>.
- Kraska, T., Beutel, A., Chi, E. H., Dean, J., and Polyzotis, N. The case for learned index structures. In *SIGMOD*. Association for Computing Machinery, 2018. URL <https://doi.org/10.1145/3183713.3196909>.
- Li, Y. and Woodruff, D. P. A tight lower bound for high frequency moment estimation with small error. In *Approximation, Randomization, and Combinatorial Optimization. Algorithms and Techniques*, pp. 623–638. Springer, 2013.
- Liu, Z., Manousis, A., Vorsanger, G., Sekar, V., and Braverman, V. One sketch to rule them all: Rethinking network flow monitoring with univmon. In *SIGCOMM*, 2016.
- Liu, Z., Ben-Basat, R., Einziger, G., Kassner, Y., Braverman, V., Friedman, R., and Sekar, V. Nitrosketch: robust and general sketch-based monitoring in software switches. In *Proceedings of the ACM Special Interest Group on Data Communication, SIGCOMM 2019, Beijing, China, August 19-23, 2019*, pp. 334–350, 2019. URL <https://doi.org/10.1145/3341302.3342076>.
- Maciá-Fernández, G., Camacho, J., Magán-Carrión, R., García-Teodoro, P., and Therón, R. UGR’16: A new dataset for the evaluation of cyclostationarity-based network idss. *Computers & Security*, 73:411 – 424, 2018. ISSN 0167-4048. doi: <https://doi.org/10.1016/j.cose.2017.11.004>. URL <http://www.sciencedirect.com/science/article/pii/S0167404817302353>.
- Manku, G. and Motwani, R. Approximate frequency counts over data streams. In *International Conference on Very Large Databases (VLDB)*, pp. 346–357, 2002.
- McGregor, A., Vorotnikova, S., and Vu, H. T. Better algorithms for counting triangles in data streams. In *Proceedings of the 35th Symposium on Principles of Database Systems (PODS)*, pp. 401–411, 2016.
- McMahan, B., Moore, E., Ramage, D., Hampson, S., and Aguera y Arcas, B. Communication-Efficient Learning of Deep Networks from Decentralized Data. In *Proceedings of the 20th International Conference on Artificial Intelligence and Statistics*, volume 54 of

*Proceedings of Machine Learning Research*. PMLR, 2017. URL <http://proceedings.mlr.press/v54/mcmahan17a.html>.

Metwally, A., Agrawal, D., and El Abbadi, A. Efficient computation of frequent and top-k elements in data streams. In *ICDT*, 2005.

Misra, J. and Gries, D. Finding repeated elements. Technical report, Cornell University, 1982.

Ohlsson, E. Sequential poisson sampling. *J. Official Statistics*, 14(2):149–162, 1998.

Paranjape, A., Benson, A. R., and Leskovec, J. Motifs in temporal networks. In *Proceedings of the Tenth ACM International Conference on Web Search and Data Mining, WSDM 17*, pp. 601–610, New York, NY, USA, 2017. Association for Computing Machinery. ISBN 9781450346757. doi: 10.1145/3018661.3018731. URL <https://doi.org/10.1145/3018661.3018731>.

Pass, G., Chowdhury, A., and Torgeson, C. A picture of search. In *Proceedings of the 1st International Conference on Scalable Information Systems, InfoScale 06*, pp. 1–es, New York, NY, USA, 2006. Association for Computing Machinery. ISBN 1595934286. doi: 10.1145/1146847.1146848. URL <https://doi.org/10.1145/1146847.1146848>.

Rosén, B. Asymptotic theory for successive sampling with varying probabilities without replacement, I. *The Annals of Mathematical Statistics*, 43(2):373–397, 1972. URL <http://www.jstor.org/stable/2239977>.

Rosén, B. Asymptotic theory for order sampling. *J. Statistical Planning and Inference*, 62(2):135–158, 1997.

## A. Datasets

For the experiments, we use the following datasets:

- AOL (Pass et al., 2006): A log of search queries collected over three months in 2006. For each query, its frequency is the number of lines in which it appeared (over the entire 92 days).
- CAIDA (CAIDA, 2016): Anonymous passive traffic traces from CAIDA’s equinix-chicago monitor. We use the data collected over one minute (2016/01/21 13:29:00 UTC), and count the number of packets for each tuple (source IP, destination IP, source port, destination port, protocol).
- Stack Overflow (SO) (Paranjape et al., 2017): A temporal graph of interactions between users on the Stack Overflow website. For each node in the graph, we consider its weighted in degree (total number of responses received that user) and its weighted out degree (total number of responses by that user).
- UGR (Maciá-Fernández et al., 2018): Real traffic information collected from the network of a Spanish ISP (for network security studies). We consider only one week of traffic (May 2016 Week 2). For a pair of source and destination IP addresses, its frequency will be the number packets sent between these two addresses (only considering flow labeled as “background”, not suspicious activity).

For the experiments with advice, we use parts of these datasets in the following way:

- AOL: We use the same predictions given in (Hsu et al., 2019), which were the result training a deep learning model on the queries from the first five days. We use the prediction to estimate frequency moments on the queries from the 51st and 81st days (after removing duplicate queries from multiple clicks on results).
- Stack Overflow: We consider two six month periods: 1/2013-6/2013 and 7/2013-12/2013. We estimate the in and out degree moments on the data from 7/2013-12/2013, where the advice for each node is its exact degree in past data (the previous six month period).

## B. Sample by Advice: Algorithm and Proofs

The pseudocode for our sampling by advice sketch and the respective estimators is provided in Algorithms 1 and 2. The sampling sketches also support merging (pseudocode not shown). Since the advice is available with each occurrence of a key, we can implement the weighted sample by advice  $a_x$  using schemes designed for aggregated data (a model where each key occurs once with its full weight) such as (Chao, 1982; Ohlsson, 1998; Duffield et al., 2007; Rosén, 1997; Cohen & Kaplan, 2008; Cohen et al., 2011). Some care (using a hash function) is needed because keys occur in multiple data elements. The pseudocode shows a ppswor implementation. The pseudocode also integrates the uniform and ppswor-by-advice samples to avoid duplication (keys that qualify for both samples are stored once).

We establish that the sampling sketch returns the exact frequencies for sampled keys:

**Lemma B.1.** *The frequency  $w_x$  of each key  $x$  in the final sample is accurate.*

*Proof.* Note that if a key enters the sample when an element

**Algorithm 1:** Sample by advice (data processing)

**Input:** A stream of updates, advice model  $a$ ,  $k_h, k_p, k_u$  (sample size for heaviest, by-advice, and uniform)

**Output:** A sampling sketch for  $f(\mathbf{w})$

**Initialization:**

**begin**

Set hash function such that  $h(x) \sim \text{Exp}(1)$  (independently for all keys) //  $\text{Exp}(1)$  for ppswor,  $U[0,1]$  for priority  
 Create empty sample  $S := (S.h, S.pu)$  //  $S.h$  stores top  $k_h$  keys by advice;  $S.pu$  is a sketch for a combined weighted-by-advice and uniform sample

**Process an update**  $(x, \Delta)$  **with prediction**  $a_x$ :

**begin**

**if**  $x \in S$  **then** //  $x$  is in (any component) of the sample  
 $w_x \leftarrow w_x + \Delta$  // Update the count of  $x$   
**Break**  
**else**  
**if**  $|S.h| < k_h$  **or**  $a_x > \min_{y \in S.h} a_y$  **then** //  $x$  has a top- $k_h$  advice value  
 insert  $(x, w_x \leftarrow \Delta)$  to  $S.h$   
**if**  $|S.h| = k_h + 1$  **then**  
 Eject  $y \leftarrow \arg \min_{z \in S.h} a_z$  from  $S.h$   
 // eject key with smallest advice in  $S.h$   
 Process  $(y, w_y)$  by the ppswor+uniform sampling sketch  $S.pu$   
**Break**  
**else**  
 process  $(x, \Delta)$  by the ppswor+uniform sampling sketch  $S.pu$

// Process update  $(y, \Delta)$  by ppswor+uniform sampling sketch  $S.pu$ :

**begin**

$r_y \leftarrow \frac{h(y)}{f(a_y)}$  // Compute by-advice seed value of key  $y$   
**if**  $r_y < \{r_z \mid z \in S.pu\}_{(k_p)}$  **then** //  $y$  has one of  $k_p$  smallest  $r_y$  and qualifies for ppswor sample  
 Insert  $(y, w_y \leftarrow \Delta)$  to  $S.pu$   
**else**  
**if**  $h(y) < \{h(z) \mid z \in S.pu\}_{(k_u)}$  **then** //  $y$  has one of  $k_u$  smallest  $h(y)$  and qualifies for uniform sample  
 Insert  $(y, w_y \leftarrow \Delta)$  to  $S.pu$   
**foreach**  $z \in S.pu \mid h(z) > \{h(z') \mid z' \in S.pu\}_{(k_u)}$  **and**  $r_z > \{r_{z'} \mid z' \in S.pu\}_{(k_p)}$  **do**  
 Eject  $z$  from  $S.pu$

**Algorithm 2:** Sample by advice (estimator computation)

**Input:** A by-advice sampling sketch for  $f(\mathbf{w})$  with parameters  $k_h, k_p, k_u$

**Output:** Sparse representation of  $f(\mathbf{w})$  and estimate of  $\|f(\mathbf{w})\|_1$

**foreach**  $x$  in  $S_h$  **do**

$\widehat{f}(w_x) \leftarrow f(w_x)$

**foreach**  $x \in S.pu$  **do** // keys stored in uniform/ppswor

samples

$\tau_p \leftarrow \{r_z \mid z \in S.pu \setminus \{x\}\}_{(k_p-1)}$  // The  $(k_p - 1)^{th}$  smallest seed of a key other than  $x$

$\tau_u \leftarrow \{h(z) \mid z \in S.pu \setminus \{x\}\}_{(k_u-1)}$  // The  $(k_u - 1)^{th}$  smallest hash value of a key other than  $x$

**if**  $h(x) < \tau_u$  **or**  $r_x < \tau_p$  **then** // key  $x$  strictly included in the uniform or by-advice samples

$p_x \leftarrow \Pr_h[h(x) < \max\{f(a_x)\tau_p, \tau_u\}]$  // For ppswor  
 $1 - e^{-\max\{\tau_u, f(a_x)\tau_p\}}$ ; For priority

$\min\{\max\{\tau_u, f(a_x)\tau_p\}, 1\}$

$\widehat{f}(w_x) \leftarrow \frac{f(w_x)}{p_x}$

Return  $\{(x, \widehat{f}(w_x))\}$  (sparse representation of  $f(\mathbf{w})$ ); The sum of  $\widehat{f}(w_x)$  as an estimate of  $\|f(\mathbf{w})\|_1$ . // For  $x$  assigned with  $\widehat{f}(w_x)$

with the key is processed and remains stored, its count is going to be accurate (we account for all the updates involving that key). Since the prediction  $a_x$  is consistent (the prediction  $a_x$  is the same in all updates involving  $x$ ), the seed of  $x$  is the same every time  $x$  appears. For  $x$  to not enter the sample on its first occurrence or to be removed at any point, there must be other  $k$  keys in the sample with seed values lower than that of  $x$ . If such keys exist,  $x$  is not in the final sample. The argument for the  $k_h$  top advice keys and for the  $k_u$  uniformly samples keys is similar.  $\square$

*Proof of Lemma 3.1.* It suffices to establish the upper bound

$$\text{Var}[\widehat{f(w_x)}] \leq f(w_x)f(\mathbf{w}) \min\left\{\frac{c_p}{k_p-2}, \frac{c_u}{k_u-2}\right\}.$$

on the variance of the estimate of each key  $x$ .

A key  $x$  with one of the top  $h$  advised frequencies has  $\text{Var}[\widehat{f(w_x)}] = 0$  and the claim trivially holds. Otherwise, recall that we assume that for some  $c_p, c_u \geq 0$ ,

$$\frac{f(w_x)}{\|f(\mathbf{w})\|_1} \leq \max\left\{c_p \frac{f(a_x)}{\|f(\mathbf{a})\|_1}, c_u \frac{1}{n}\right\}.$$

The variance is  $\text{Var}[\widehat{f(w_x)}] = f(w_x)^2 \mathbf{E}_{p_x}[(1/p_x - 1)]$ , where  $p_x$  is as computed by the algorithm. Now note that  $p_x = \max\{p'_x, p''_x\}$ , where  $p'_x$  is the probability  $x$  is included in a ppswor-by-advice sample of size  $k_p$  and  $p''_x$  is the probability it is included in a uniform sample of size  $k_u$ . We obtain that the variance is the minimum of the variance in these two scenarios.

The variance from a uniform sample of size  $k_u$  is bounded by  $\frac{1}{k_u-2}nf(w_x)^2$ . If  $\frac{f(w_x)}{\|f(\mathbf{w})\|_1} \leq c_u \frac{1}{n}$  we substitute  $f(w_x) \leq \|f(\mathbf{w})\|_1 c_u \frac{1}{n}$  and obtain

$$\begin{aligned} \frac{1}{k_u-2}nf(w_x)^2 &\leq \frac{1}{k_u-2}nf(w_x)\|f(\mathbf{w})\|_1 c_u \frac{1}{n} \\ &= \frac{c_u}{k_u-2}f(w_x)\|f(\mathbf{w})\|_1. \end{aligned}$$

The variance from a weighted-by-advice sample of size  $k_p$  is bounded by

$$\frac{1}{k_p-2}f(a_x)\|f(\mathbf{a})\|_1 \frac{f(w_x)^2}{f(a_x)^2} = \frac{1}{k_p-2}\|f(\mathbf{a})\|_1 \frac{f(w_x)^2}{f(a_x)}.$$

If  $\frac{f(w_x)}{\|f(\mathbf{w})\|_1} \leq c_p \frac{f(a_x)}{\|f(\mathbf{a})\|_1}$  we similarly substitute and obtain that the variance is at most

$$\frac{c_p}{k_p-2}f(w_x)\|f(\mathbf{w})\|_1.$$

□

### C. Sample by Advice: Experiments

Results for performance of sampling by advice on additional datasets are provided in Figure 5.

#### Error estimation

The estimators we use are unbiased and we consider the Normalized Mean Squared Error (NRMSE) which for unbiased estimators is the same as the Coefficient of Variation (CV), the ratio of the standard deviation to the mean:

$$\frac{\text{Var}[\|\widehat{f(\mathbf{w})}\|_1]^{1/2}}{\|f(\mathbf{w})\|_1}.$$

A simple way to estimate the variance is to use the empirical squared error over a set of runs: We take the average of  $(\|\widehat{f(\mathbf{w})}\|_1 - \|f(\mathbf{w})\|_1)^2$  over runs and apply a square root. We found that 50 runs were not sufficient for an accurate estimate with our sample-by-advice methods. This is because of keys that had relatively high frequency and low advice, which resulted in low inclusion probability and high contribution to the variance. Samples that included these keys had higher estimates of the statistics than the bulk of other samples and often significant increase was attributed to one or two keys. This could be remedied by significantly increasing the number of runs we average over. We instead opted to use different and more accurate estimators for the by-advice and (for consistency) for the baseline with and without-replacement schemes.

For with-replacement schemes we computed the exact variance (and hence the NRMSE) as follows: The inclusion probability of a key  $x$  in a sample of size  $k$  is

$$p'_x := 1 - (1 - p_x)^k$$

where  $p_x$  is the probability to be selected in one sampling step. That is,  $p_x = w_x/\|\mathbf{w}\|_1$  for with-replacement  $\ell_1$  sampling and  $p_x = w_x^2/\|\mathbf{w}\|_2^2$  for with-replacement  $\ell_2$  sampling. The per-key variance of our estimator for key  $x$  is  $(1/p'_x - 1)(f(w_x))^2$ . Since estimates for different keys have 0 correlations, the variance of our estimate of the statistics  $\|f(\mathbf{w})\|_1$  is

$$\text{Var}[\|\widehat{f(\mathbf{w})}\|_1] := \sum_x (1/p'_x - 1)f(w_x)^2.$$

For the without-replacement schemes (the by-advice sampling and the without-replacement reference methods) we apply a more accurate estimator over the same set of 50 runs. For each "run" and each key  $x$  (sampled or not) we consider all the possible samples where the randomization of all keys  $y \neq x$  is as in the "run." These include samples that include and do not include  $x$ . We then compute the inclusion probability  $p'_x$  of key  $x$  under these conditions. The contribution to the variance due to this set of runs is  $(1/p'_x - 1)f(w_x)^2$ . We sum the estimates  $(1/p'_x - 1)f(w_x)^2$  over all keys and take the average of the sums over runs as our estimate of the variance.

For the pure bottom- $k$  (the ppswor by  $w^q$  for  $q = 1, 2$ ) recall that we compute random *seed* values to keys of the form  $r_x/w_x^q$ , where  $r_x \sim D$  are independent. The sample is the  $k$  keys with lowest seed values. The inclusion probability is computed with respect to a threshold that is defined to be the  $k$ th smallest "seed" value of other keys  $\tau_x \leftarrow \{\text{seed}(y) | y \neq x\}_{(k)}$ . For keys in the sample, this is the  $(k+1)^{\text{th}}$  smallest seed overall. For keys not in the sample the applicable  $\tau_x$  is the  $k^{\text{th}}$  smallest overall. We get  $p'_x = \Pr_{r \sim D}[r_x/w_x^q \leq \tau_x]$  (a different  $p'_x$  is obtained in each run).

The calculation for the by-advice+uniform sampling is as in the estimator in Algorithm 1, except that it is computed for all keys  $x$ .

### D. Combinations of Frequencies and Functions: Proofs and Details

#### D.1. Near-uniform Frequency Distributions

We showed that frequency distributions with heavy hitters are easy for high moments and moreover, the validity of the result can be certified. Interestingly, the other extreme of near-uniform distributions (where  $w_1/w_n$  is bounded) are also easy. But unfortunately, unlike the case with heavy hitters, there is no "certificate" to the validity of the emulation.

**Lemma D.1.** *Let  $w$  be a frequency distribution with support size  $n$ . Then the overhead of using  $\ell_1$  or  $\ell_0$  sampling to emulate  $\ell_p$  sampling is at most  $\left(\frac{w_1}{w_n}\right)^p$ .*

*Proof.* We use Lemma 4.1 and the full form of the over-

head bound (5) and lower bound the denominator  $\left\| \frac{\mathbf{w}}{w_1} \right\|_p^p$ .

Note that for any  $\mathbf{w}$  with support size  $n$ ,  $\left\| \frac{\mathbf{w}}{w_1} \right\|_1^1 \leq n$  and

$\left\| \frac{\mathbf{w}}{w_1} \right\|_0 = n$ . Now,

$$\left\| \frac{\mathbf{w}}{w_1} \right\|_p^p = \sum_{i=1}^n \left( \frac{w_i}{w_1} \right)^p \geq \sum_{i=1}^n \left( \frac{w_n}{w_1} \right)^p = n \cdot \left( \frac{w_n}{w_1} \right)^p.$$

The overhead for  $q \in \{0, 1\}$  is then

$$\frac{\left\| \frac{\mathbf{w}}{w_1} \right\|_q^q}{\left\| \frac{\mathbf{w}}{w_1} \right\|_p^p} \leq \frac{n}{n \cdot \left( \frac{w_n}{w_1} \right)^p} = \left( \frac{w_1}{w_n} \right)^p.$$

□

## E. Additional Experiments

Estimates of the distribution of frequencies for all datasets are reported in Figure 6. The estimates are from ppswor ( $\ell_1$  without replacement) and  $\ell_2$  (with replacement) samples of sizes  $k = 32$  and  $k = 1024$ . For each key  $x$  in the sample, we estimate its rank in the data set (the number of keys  $y$  with frequency  $w_y \geq w_x$ ). The estimate is computed using the threshold function  $f(w) := I_{w \geq w_x}$ . The pairs of frequency and estimated rank are then plotted. The figures also provide the exact frequency distribution.

Additional results on estimation of moments from ppswor ( $\ell_1$  without replacement) and  $\ell_2$  samples (with and without replacement) are reported in Figure 7. As suggested by our analysis, the estimates on all datasets are surprisingly accurate even with respect to the "benchmark" upper bound (for weighted with-replacement sampling tailored to the moment we are estimating). The figures also show the advantage of "without replacement" sampling on these skewed datasets.



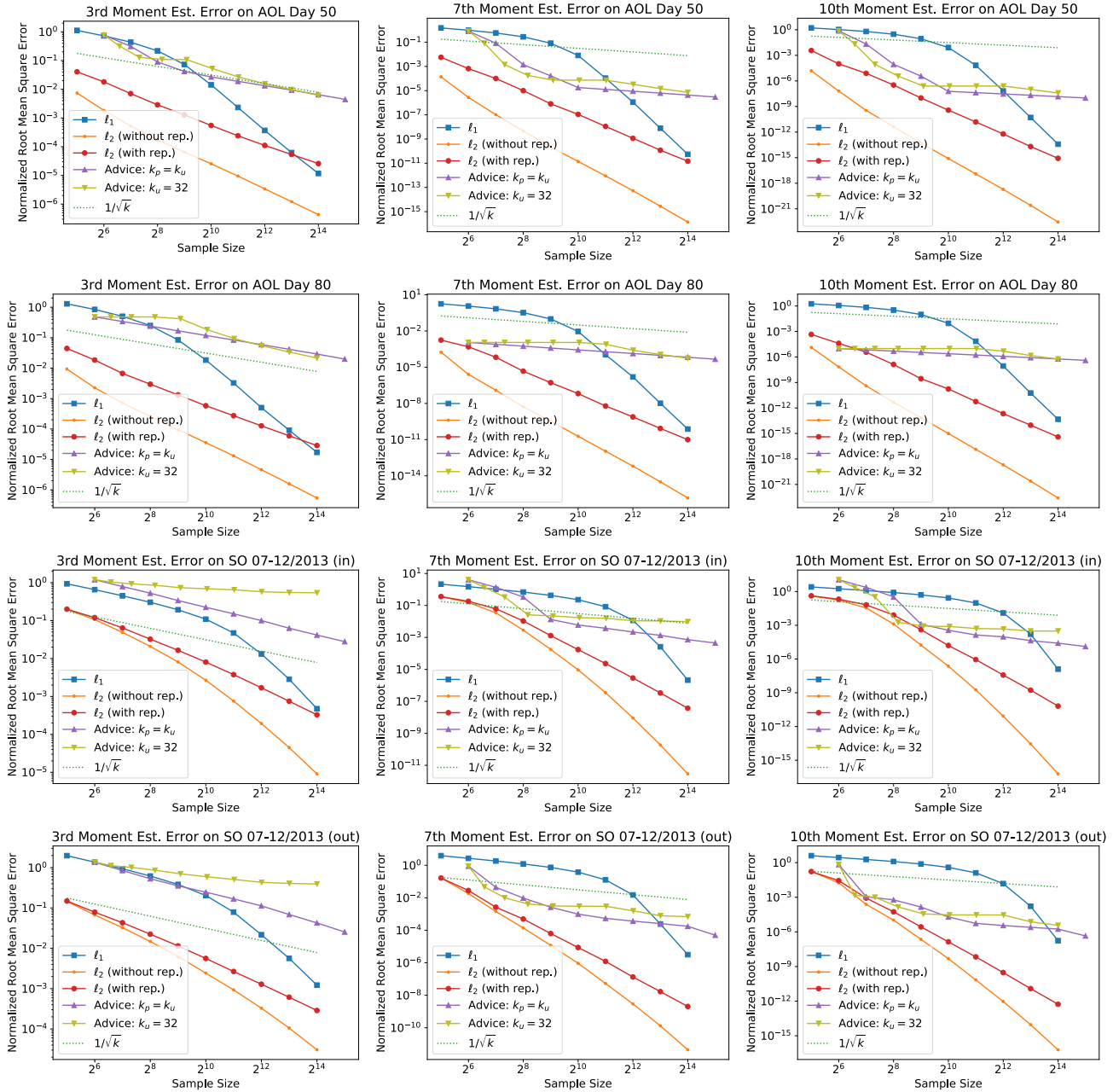


Figure 5. NRMSE for estimating 3rd, 7th, and 10th moments on AOL dataset (days 50 and 80 with learned advice from (Hsu et al., 2019)) and on Stack Overflow dataset (outgoing and incoming edges), based on past frequencies.

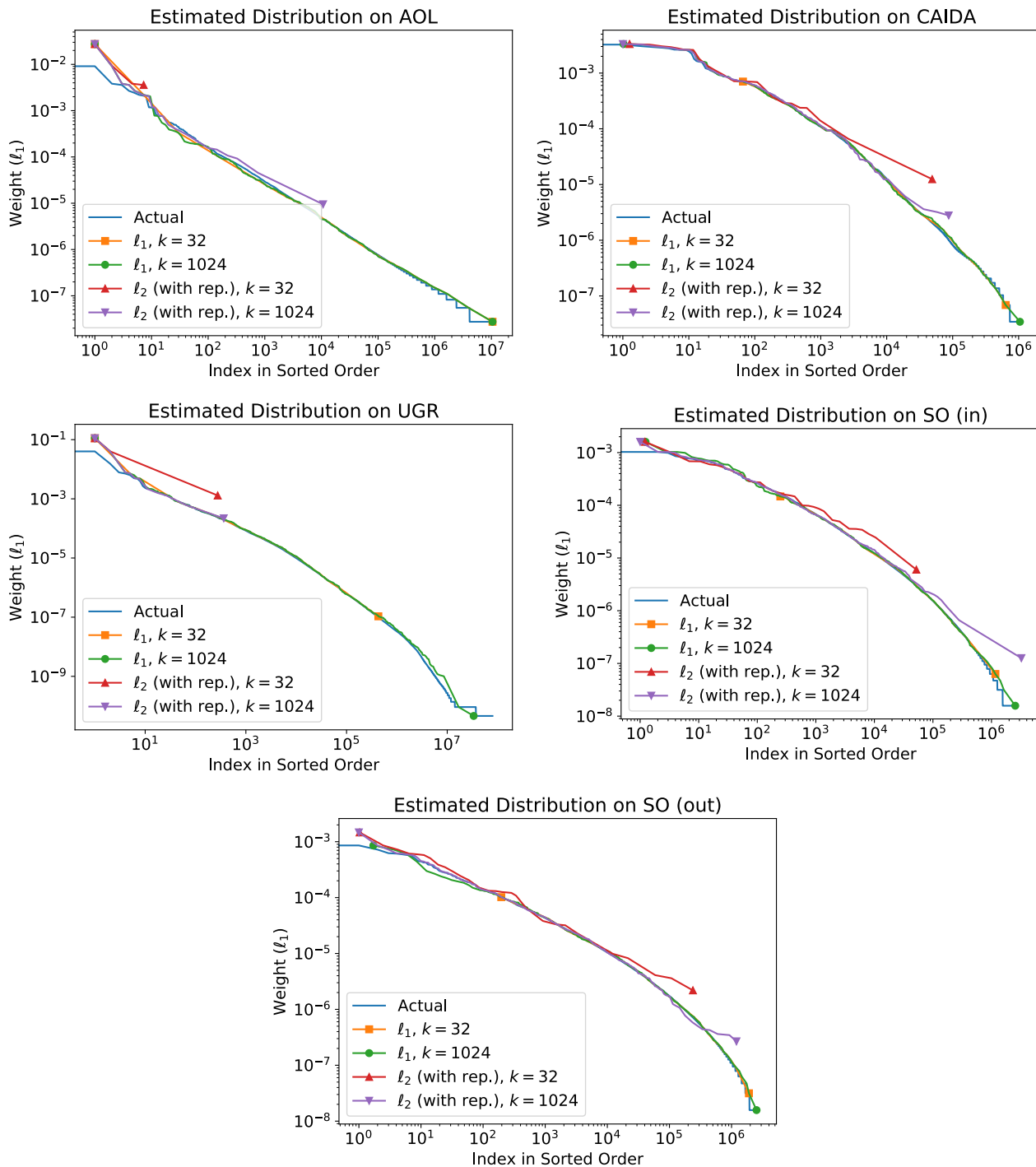


Figure 6. Actual and estimated distribution of frequency by rank. Estimates for ppswor and with-replacement  $\ell_2$  sampling and sample sizes  $k = 32, 1024$ .

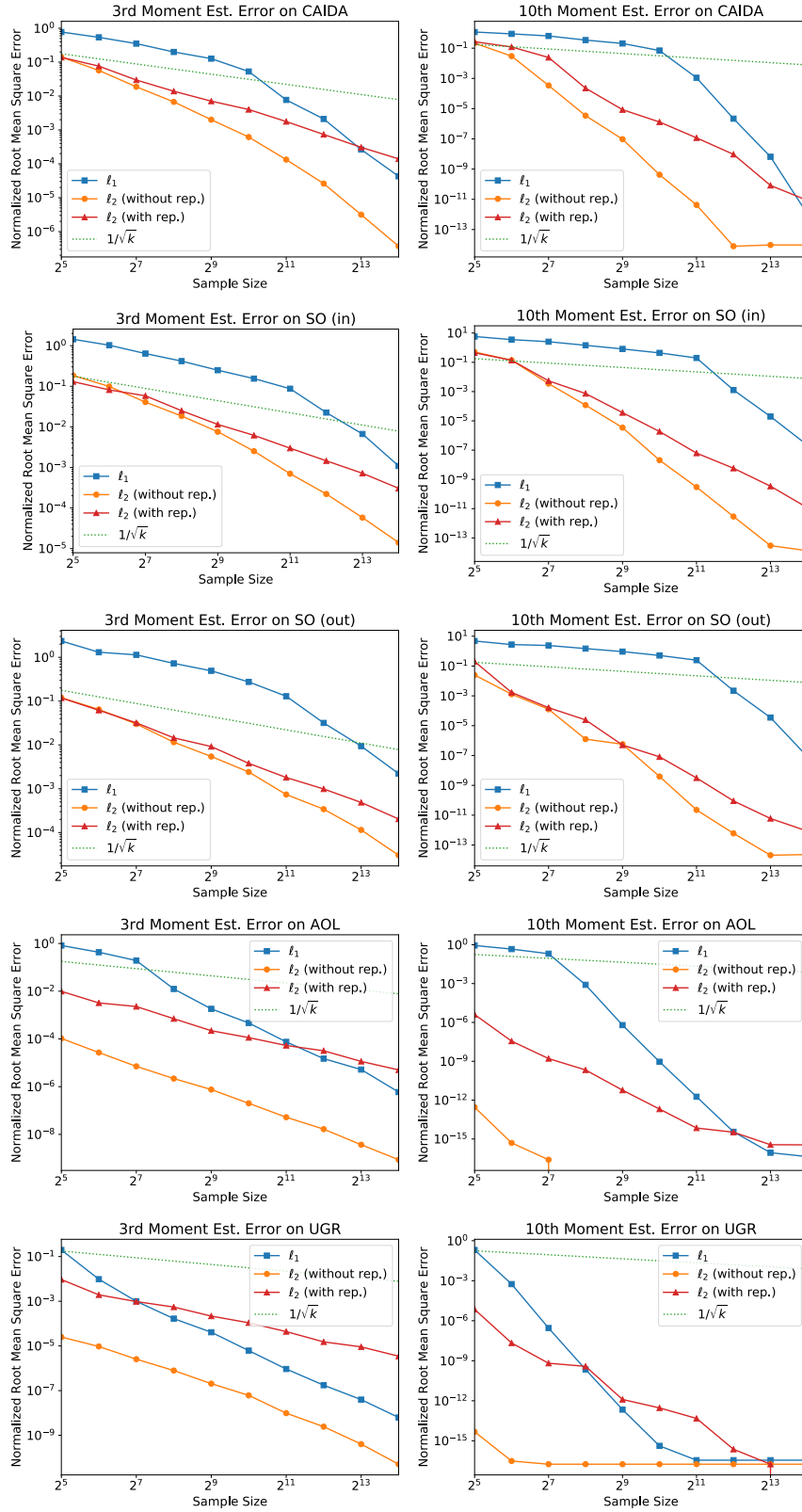


Figure 7. Estimating 3rd and 10th moments on various datasets from ppswor ( $\ell_1$  without replacement) and  $\ell_2$  samples (with and without replacement). The error is averaged over 50 repetition.