


Efficient Differentially Private F_0 Linear Sketching

Rasmus Pagh 

IT University of Copenhagen
 BARC
 pagh@itu.dk

Nina Mesing Stausholm 

IT University of Copenhagen
 BARC
 nimm@itu.dk

Abstract

A powerful feature of *linear sketches* is that from sketches of two data vectors, one can compute the sketch of the difference between the vectors. This allows us to answer fine-grained questions about the difference between two data sets. In this work we consider how to construct sketches for weighted F_0 , i.e., the summed weights of the elements in the data set, that are small, differentially private, and computationally efficient. Let a weight vector $w \in (0, 1]^u$ be given. For $x \in \{0, 1\}^u$ we are interested in estimating $\|x \circ w\|_1$ where \circ is the Hadamard product (entrywise product).

Building on a technique of Kushilevitz et al. (STOC 1998), we introduce a sketch (depending on w) that is linear over $\text{GF}(2)$, mapping a vector $x \in \{0, 1\}^u$ to $Hx \in \{0, 1\}^\tau$ for a matrix H sampled from a suitable distribution \mathcal{H} . Differential privacy is achieved by using *randomized response*, flipping each bit of Hx with probability $p < 1/2$. That is, for a vector $\varphi \in \{0, 1\}^\tau$ where $\Pr[(\varphi)_j = 1] = p$ independently for each entry j , we consider the *noisy sketch* $Hx + \varphi$, where the addition of noise happens over $\text{GF}(2)$. We show that for every choice of $0 < \beta < 1$ and $\varepsilon = O(1)$ there exists $p < 1/2$ and a distribution \mathcal{H} of linear sketches of size $\tau = O(\log^2(u)\varepsilon^{-2}\beta^{-2})$ such that:

1. For random $H \sim \mathcal{H}$ and noise vector φ , given $Hx + \varphi$ we can compute an estimate of $\|x \circ w\|_1$ that is accurate within a factor $1 \pm \beta$, plus additive error $O(\log(u)\varepsilon^{-2}\beta^{-2})$, w. p. $1 - u^{-1}$, and
2. For every $H \sim \mathcal{H}$, $Hx + \varphi$ is ε -differentially private over the randomness in φ .

The special case $w = (1, \dots, 1)$ is *unweighted* F_0 . Previously, Mir et al. (PODS 2011) and Kenthapadi et al. (J. Priv. Confidentiality 2013) had described a differentially private way of sketching unweighted F_0 , but the algorithms for calibrating noise to their sketches are not computationally efficient, either using quasipolynomial time in the sketch size or superlinear time in the universe size u .

For fixed ε the size of our sketch is polynomially related to the lower bound of $\Omega(\log(u)\beta^{-2})$ bits by Jayram & Woodruff (Trans. Algorithms 2013). The additive error is comparable to the bound of $\Omega(1/\varepsilon)$ of Hardt & Talwar (STOC 2010). An application of our sketch is that two sketches can be added to form a noisy sketch of the form $H(x_1 + x_2) + (\varphi_1 + \varphi_2)$, which allows us to estimate $\|(x_1 + x_2) \circ w\|_1$. Since addition is over $\text{GF}(2)$, this is the weight of the symmetric difference of the vectors x_1 and x_2 . Recent work has shown how to privately and efficiently compute an estimate for the symmetric difference size of two sets using (non-linear) sketches such as FM-sketches and Bloom Filters, but these methods have an error bound no better than $O(\sqrt{\bar{m}})$, where \bar{m} is an upper bound on $\|x_1\|_0$ and $\|x_2\|_0$. In particular, our result improves previous work when $\beta = o(1/\sqrt{\bar{m}})$ and $\log(u)/\varepsilon = \bar{m}^{o(1)}$.

In conclusion our results both improve the efficiency of existing methods for unweighted F_0 estimation and extend to a weighted generalization. We also give a distributed streaming implementation for estimating the size of the union between two input streams.

2012 ACM Subject Classification Security and privacy \rightarrow Formal methods and theory of security

Keywords and phrases Differential Privacy, Linear Sketches, Weighted F_0 Estimation

Funding This work was supported by Investigator Grant 16582, Basic Algorithms Research Copenhagen (BARC), from the VILLUM Foundation.

1 Introduction

Estimating the number of distinct values in a set (its *cardinality*), without explicitly enumerating the set, is a classical and important problem in data management. Sampling-based methods [23] can in many cases be improved by using algorithms designed with data streams in mind [26]. Streaming algorithms based on *linear sketches* can also be used to estimate changes as a data set evolves [28] and for approximate query processing in distributed settings [3, 13]. As our first motivating example consider the following SQL query:

```
SELECT P.name
FROM EMPLOYEES E, HOSPITALIZATION H
WHERE E.salary > 100000 AND E.name = H.name AND H.year = 2020
```

The size (in bytes) of the query result is a sum weighted by string length over the names that appear in subsets of two relations. That is, estimating the size of the join result is about estimating the *weighted* size of a set intersection.

In recent years, *privacy* of database records has become increasingly important when releasing aggregates from a database. In the example above, the information that a tuple with a particular person exists (and satisfies a certain predicate) can potentially be sensitive. If the database is distributed, with relations on different servers that are not allowed to expose sensitive information, it is not trivial how to even estimate the join size.

The notion of *differential privacy* [16] has emerged as the leading approach to providing rigorous privacy guarantees. It is known that differential privacy comes with pitfalls [29], but work in the database community has led to privacy-preserving database systems supporting (limited) SQL, see e.g. [33, 47] and their references. A challenge in such systems is that the set of queries is often not known ahead of time, so *budgeting* the disclosure of detailed information is highly nontrivial. An attractive approach to achieving privacy even when faced with unknown queries is to release a summary, or sketch, of the data set from which approximate answers to queries can be computed (as a side effect this also eliminates the need for interaction). In this paper we consider private linear sketches for the problem of cardinality estimation.

Example. Suppose that the company Acme Corporation runs an employee satisfaction survey once a year. Management at Acme Corporation made some drastic changes over the past year, and they wish to analyze the impact of these changes on the employees' satisfaction. For a specific improvement, every employee is given a value between 0 and 1, indicating how closely related that improvement is to the employee's work life. A survey for each improvement is run by a consultant who delivers a summary of the results to the management at Acme Corporation. The consultant ensures that the summary is private, so individual employees cannot be identified from the summary. The management at Acme Corporation can combine the summary from last year's survey with the summary from this year's survey to estimate the change in satisfaction over the past year, where the vote of an employee is weighted by the value that employee was given. We note that the summaries should be generated in the same way, but the choice of consultant may change from year to year.

More formally, we consider two players that hold sets A and B from a universe $U = \{1, \dots, u\}$, respectively. For every element $j \in U$ let a fixed, public weight, $w_j \in (0, 1]$ be given and for input set $A \subseteq U$ consider the corresponding weight vector $(w_A)_j = w_j \cdot \mathbf{1}[j \in A]$. The goal is to estimate the weight of the symmetric difference $\|w_{A \Delta B}\|_1$, in a differentially

private manner. We refer the reader to Section 3.2 for the basics of differential privacy. We may think of the sets as two lists of employees. Given input sets A and B , the two players each compute a *linear sketch* of their own set and add noise to obtain privacy as described in Section 4. These noisy sketches can be thought of as the summaries.

For input sets A and B , we note that if we, along with the estimate of the weight of the symmetric difference, have estimates of $\|w_A\|_1$ and $\|w_B\|_1$, then we can also estimate $\|w_{A \cup B}\|_1$, $\|w_{A \cap B}\|_1$, $\|w_{A \setminus B}\|_1$ and $\|w_{B \setminus A}\|_1$ as argued in Section 4.3. To make this possible, each player also outputs a differentially private version of their set weight. We remark that if all weights $w_j = 1$, then the problem reduces to estimating the set *size*, a problem often referred to as F_0 .

We define and construct a noisy linear sketch over $\text{GF}(2)$, the field of size 2, with the following properties:

- ε -differentially private
- Computationally efficient
- Allows estimating the weight of the symmetric difference with small relative error
- Space usage is polynomially related to the lower bound (for fixed ε)

Previously known results satisfy at most 3 of these properties, see Figure 1 for an overview. We discuss previous work further in Section 2. Our sketch can be computed and stored for future use, meaning that two players do not have to be active simultaneously but can compute and publish their sketches when they are ready. A self-contained description of our linear sketch can be found in Section 4. Readers familiar with the sketching literature will realize that our sketch combines a method of Kushilevitz, Ostrovsky, and Rabani [30] with a standard hashing-based subsampling technique (see, e.g., [48]), and we use a Randomized Response Technique [46] with noise parameter $p(\varepsilon)$, to get ε -differential privacy. Hence, refer to our sketch as the *KOR sketch* and to its noisy counterpart as a *noisy KOR sketch*. We note that a related, but non-linear and non-private, sketch has previously been used for estimating size of symmetric difference [39]. From now on we leave out ε in the noise parameter and write simply p . We show that the KOR sketch is sufficiently robust to noise to allow precise estimation after adding noise, thus allowing pure differential privacy.

We next give an overview of our techniques, discussed in depth in Section 4. Let $U = \{1, \dots, u\}$ be the universe from which the input sets are taken. Privacy parameter ε and accuracy parameter β are given, and a sketch size τ is determined by these parameters. We show in Section 5.2 that we can construct an ε -differentially private sketch from which we can compute a $(1 + \beta)$ -approximation for the weight of the symmetric difference with high probability.

Randomized response [46] is applied to the entire sketch Hx , meaning that each entry of the sketch is flipped with probability $p < 1/2$. We show in Section 5.1 how to choose p as a function of ε to ensure ε -differential privacy for the sketch. Let $x \circ w$ denote the Hadamard product. Our main theorem is:

► **Theorem 1** (Noisy KOR sketch). *Let $w \in (0, 1]^u$ be given. For every choice of $0 < \beta < 1$ and $\varepsilon = O(1)$ there exists a distribution \mathcal{H} over $\text{GF}(2)$ -linear sketches mapping a vector $x \in \{0, 1\}^u$ to $\{0, 1\}^\tau$, where $\tau = O(\log^2(u)\varepsilon^{-2}\beta^{-2})$, and a distribution \mathcal{N}_ε over noise vectors such that:*

1. *For $H \sim \mathcal{H}$ and $\varphi \sim \mathcal{N}_\varepsilon$, given $Hx + \varphi$ we can compute, in time $O(\tau)$, an estimate \hat{w} of $\|x \circ w\|_1$ that with probability $1 - 1/u$ satisfies $|\hat{w} - \|x \circ w\|_1| < \beta \|x \circ w\|_1 + O(\log(u)\varepsilon^{-2}\beta^{-2})$.*
2. *For every H in the support of \mathcal{H} , $Hx + \varphi$ is ε -differentially private over the choice of $\varphi \sim \mathcal{N}_\varepsilon$, and can be computed in time $O(\|x\|_0 \log(u) + \tau)$, including time for sampling φ .*

Reference	Diff. privacy	Additive error	Relative error	Initial. time	Space usage
Hardt and Talwar [24]	ε	$\Omega(1/\varepsilon)$	–	–	–
McGregor et al. [31]	ε	$\tilde{\Omega}(\sqrt{m}/e^\varepsilon)$	–	–	–
Jayram and Woodruff [25]	–	–	$1 + \beta$	–	$\tilde{\Omega}(1/\beta^2)$
Kane et al. [26]	–	$\tilde{O}(1)$	$1 + \beta$	$O(1)$	$\tilde{O}(1/\beta^2)$
Mir et al. [36]	ε	$\tilde{O}(m^{1-\Omega(1)}/\varepsilon^{O(1)})$	$1 + \beta$	$\exp((\varepsilon\beta)^{-O(1)})$	$\tilde{O}((\varepsilon\beta)^{-O(1)})$
Kenthapadi et al. [27]	(ε, δ)	$\tilde{O}(\sqrt{m}/\varepsilon)$	$1 + \beta$	$\tilde{\Omega}(u)$	$\tilde{O}(1/\beta^2)^*$
Stanojevic et al. [42]	ε	$\tilde{O}(\sqrt{ A \cup B }/\varepsilon^2)$	–	$\Omega(A + B)$	$\Omega(A + B)$
This paper	ε	$\tilde{O}(m^{2/3}/\varepsilon^{2/3})$	$1 + \beta$	$\tilde{O}(\varepsilon^{-2}\beta^{-2})$	$\tilde{O}(\varepsilon^{-2}\beta^{-2})$

■ **Figure 1** Selected lower bounds (top part) and upper bounds (bottom part) for estimating the (unweighted) size of the symmetric difference $m = |A \Delta B|$ from small sketches of sets $A, B \subseteq \{1, \dots, u\}$. Bounds stated as \tilde{O} and $\tilde{\Omega}$ are simplified by suppressing multiplicative factors polynomial in $\log(1/\varepsilon)$, $\log(1/\beta)$, $\log(1/\delta)$, and $\log u$. The non-private bounds in [25, 26] improve previous results by an $\tilde{O}(1)$ factor, we refer to their references for details. * The space usage of [27] is measured in terms of real numbers; it is unclear how much space a private, discrete implementation would need.

The assumption that $\varepsilon = O(1)$ is not essential, and is only made to simplify our bounds (which do not improve for privacy parameter $\varepsilon = \omega(1)$). Without loss of generality we can assume that parameter β is such that the error is dominated by $\beta \|x \circ w\|_1$, because reducing β further cannot reduce error by more than a factor 2. In the unweighted case, setting $\beta = \sqrt[3]{\log(u)/(\varepsilon^2 m)}$ to balance relative and additive error we get error $\tilde{O}(m^{2/3}/\varepsilon^{2/3})$, where the \tilde{O} notations suppresses a polylogarithmic factor. This is polynomially related to known lower bounds described in section 2.3.

Applications

Suppose that Alice holds set A with corresponding characteristic vector $x_A \in \{0, 1\}^u$ and Bob holds set B with characteristic vector $x_B \in \{0, 1\}^u$. They jointly sample $H \sim \mathcal{H}$ and privately sample $\varphi_A, \varphi_B \sim \mathcal{N}_\varepsilon$ according to Theorem 1. Then $Hx_A + \varphi_A$ and $Hx_B + \varphi_B$ are ε -differentially private. Furthermore, $(Hx_A + \varphi_A) + (Hx_B + \varphi_B) = (Hx_A + Hx_B) + (\varphi_A + \varphi_B)$, and we show in Section 4.3 that $\varphi_A + \varphi_B \sim \mathcal{N}_{\varepsilon'}$ with $\varepsilon' = \varepsilon^2/(2 + 2\varepsilon)$. In Section 5.2 we use this in conjunction with Theorem 1 to establish:

► **Corollary 2.** *For accuracy parameter $\beta > 0$, consider an ε -differentially private noisy KOR sketch for a set A and an ε -differentially private noisy KOR sketch for a set B , based on the same linear sketch $H \sim \mathcal{H}$, sampled independently of A and B . We can compute an approximation $\hat{\Delta}$ of the weight of the symmetric difference, such that with probability $1 - 1/u$:*

$$\|w_{A \Delta B}\|_1 - \hat{\Delta} < \beta \|w_{A \Delta B}\|_1 + \text{poly}(1/\varepsilon, 1/\beta, \log u) .$$

In the special case where all weights w_j are 1, this reduces to estimating the *size* of the symmetric difference $A \Delta B$.

In Section 6 we describe how to modify our sketch to apply in a streaming setting. In this case, we estimate the size of the union of the input streams rather than the size of the symmetric difference when merging two sketches.

2 Related Work

In the absence of privacy constraints, seminal estimators for (unweighted) set cardinality that support merging sketches (to produce a sketch of the union) are HyperLogLog [20], FM-sketches [21], and bottom- k (aka. k -minimum values) sketches [6]. Progress on making these estimators private for set operations include [43] (using FM-sketches) and [41], which builds a private cardinality estimator to estimate set intersection size using the bottom- k sketch. We note that these sketches do not achieve differential privacy, but are aimed at a weaker notion of privacy. Specifically, they offer a one-sided guarantee that may reveal that an individual element is *not* present in the dataset. To our best knowledge, a private version of HyperLogLog with provable bounds on accuracy has not been described in the literature.

The *weighted* version of cardinality estimation has been less studied. For (scaled) integer weights in $[W]$ there is a simple reduction that inserts element i with weight w_i by inserting the tuples $(i, 1), \dots, (i, w_i)$ into a standard cardinality estimator on the domain $U \times [W]$, but this makes the obtained bounds depend on the number W of possible weights. Cohen et al. [12] showed that the class of cardinality estimators that rely on extreme order statistics (for example HyperLogLog) can be efficiently extended to the weighted setting, even for real-numbered weights.

Note that the weighted F_0 estimation problem is different from F_1 and L_1 estimation in the context of set operations, for example, the union of two identical sets will have the same weighted F_0 , whereas summing two identical vectors will produce a vector with twice the L_1 norm. In the rest of this section we focus on the standard, unweighted setting.

2.1 Differentially private cardinality estimators

Already the seminal paper on pan-privacy [17] discusses differentially private streaming algorithms for F_0 on insertion-only streams. Their sketch is not linear and does not allow deletions or subtraction of sketches. It is not clear if the sketch can be merged to produce a sketch for the union. Recent work by von Voigt et al. [45] has shown how to estimate the cardinality of a set using less space in a differentially private manner using FM-sketches, using the Probabilistic Counting with Stochastic Averaging (PCSA) technique [21]. These sketches can be merged to obtain a sketch for the union of the input set with a slightly higher level of noise. Privacy is achieved by randomly adding ones to the sketch and by only sketching a sample of the input dataset.

Bloom Filters have been studied extensively to obtain cardinality estimators under set operations (already implicit in [17]). Alaggar et al. [2] estimated set intersection size by combining a technique for computing similarity between sets, represented by Bloom filters in a differentially private manner, named BLIP (BLoom-then-FLIP) filters [1] with a technique for approximating set intersection of two sets based on their Bloom Filter representation [9]. We note that [1] achieves privacy by flipping each bit of the Bloom filter with a certain probability, much like the technique we use to get privacy of our sketch. Stanojevic et al. [42] show how to estimate set intersection, union and symmetric difference for two sets by computing an estimate for the size of the union, and combined with the size of each set, they show how to compute an estimate for the size of the intersection and the symmetric difference. They achieve privacy by flipping each bit with some probability, like in [1]. Also, RAPPOR [19] uses Bloom Filters with a Randomized Response technique to collect data from users in a differentially private way but is mainly aimed at computing heavy hitters.

Though a bound on the expected worst-case error of privately estimating the size of a symmetric difference $|A \Delta B|$ (as in Corollary 2) is not stated in any of these papers, an

upper bound of $O(\sqrt{\bar{m}})$, where \bar{m} is an upper bound on the size of the sets, follows from the discussion in [42] (for fixed ε). It seems that this magnitude of error is inherent to approaches using Bloom filters since it arises by balancing the error related to the noise and the error related to hash collisions in the Bloom filter. An advantage and special case of our noisy KOR sketch is that it can be used to directly estimate the size of the symmetric difference, and so the error will depend only on the size of the symmetric difference. It seems that with non-linear sketches it would be necessary to first estimate the size of the union and combine this with the size of each input set as exhibited in, for example, [42]. Hence, the error would depend on the size of the union of the input sets.

2.2 Differentially private sketches

Closely related to our work is the differentially private Johnson-Lindenstrauss (JL) sketch by Kenthapadi et al. [27], in which the technique of adding noise to the sketch is also applied. Kenthapadi et al. add Gaussian noise, so to store and maintain a sketched vector, some kind of discretization would be needed (not discussed in their paper). Discretizing a real-valued private mechanism is non-trivial: Without sufficient care, one might lose privacy due to rounding in an implementation, as argued by Mironov [37]. Even if a suitable discretization of the mechanism in [27] would be possible (see [10] for a general discussion), it has several drawbacks compared to our method:

- It only achieves *approximate* differential privacy as opposed to the pure differential privacy of the noisy KOR sketch.
- The time needed to update the sketch when a set element is inserted or removed is not constant (in the main method described it is linear in the sketch size).
- The time needed to initialize the sketch is linear in the size of the sketch matrix, which has u columns, because the noise needs to be calibrated to the sensitivity of the JL sketch matrix, which requires linear time in the size of the sketch matrix. Alternatively, which is the suggestion in Kenthapadi et al., the sketch matrix is assumed to have low sensitivity and noise is calibrated to this sensitivity. If a sketch matrix with a large entry is randomly chosen, the sensitivity of the sketch matrix is large, in which case the noise does not ensure privacy. So with a small probability, privacy is not preserved.

Another closely related work is the paper of Mir et al. [36], which also adds a noise vector after computing standard linear sketches for F_0 estimation to make the sketch differentially private. They further initialize their sketches with random noise vectors to also get pan-privacy. The error bound obtained is similar to ours, and the sketch has a discrete representation, but their method is inferior in terms of time complexity. This is because they rely on the *exponential mechanism* [32], which is not computationally efficient. (Note that a preprint of the paper of Mir et al. [35] presented a computationally more efficient method. However, the sensitivity analysis in that paper has an error [40] that was corrected in the slower method published in [36].)

Our method is more computationally efficient and arguably simpler than the methods of [27, 36]. Our linear sketch is not a replacement for these sketches, though, since our sketch is over $\text{GF}(2)$ rather than the reals (or integers).

2.3 Lower bounds.

Jayram and Woodruff [25] show that, even with no privacy guarantee, to obtain error probability $1/u$ we need a sketch of $\Omega(\log(u)\beta^{-2})$ bits to estimate F_0 with relative error $1 \pm \beta$. It is easy to extend this lower bound to our setting, in which an additive error of c is

allowed: Simply insert each item c times, to increase the size of the set so that the additive error is negligible. Formally this requires us to extend the universe to $U \times \{1, \dots, c\}$, such that the lower bound in terms of the original universe size becomes $\Omega(\log(u/c)\beta^{-2})$. (The reason why we do not use this reduction to eliminate the additive error in our upper bound is that the reduction increases the sensitivity of updates, destroying the differential privacy properties.)

Hardt and Talwar [24] show that an ε -differentially private sketch for F_0 must have additive error $\Omega(1/\varepsilon)$, which is comparable (up to polynomial and logarithmic factors) to the additive error we achieve.

Desfontaines et al. [14] show that it is not possible to preserve privacy in accurate cardinality estimators if we can merge several sketches without loss in accuracy. Our sketch will have an increase in noise when merging sketches, and thus does not satisfy the requirement for cardinality estimators formulated in [14].

McGregor et al. [31] showed that in order to estimate the size of the intersection of two sets A and B , based on differentially private sketches of A and B , an additive error of $\Omega(\sqrt{u}/e^\varepsilon)$ is needed in the worst case when A and B are arbitrary subsets of $[u]$. The lower bound holds even in an interactive setting where Alice (holding A) and Bob (holding B) can communicate, and we require that the communication transcript is differentially private. The hard input distribution uses sets with symmetric difference of size $\Theta(u)$ with high probability. Since $|A \cap B| = (|A| + |B| - |A \Delta B|)/2$, estimating the intersection size is no more difficult (up to constant factors in error) than estimating $|A|$, $|B|$, and $|A \Delta B|$. We can estimate $|A|$ and $|B|$ with error $O(1/\varepsilon)$ under differential privacy, so it follows that estimating $|A \Delta B|$ under differential privacy requires error $\Omega(\sqrt{u}/e^\varepsilon)$. For a contrasting upper bound, [44, 38] suggest an algorithm estimating two-party set intersection size up to an additive error of $O(\sqrt{u}/\varepsilon)$ with high probability. A lower bound in terms of the size m of the symmetric difference follows by setting $u = m$.

2.4 Noisy sketching.

In addition to the paper of Mir et al. [36], there is some previous work on sketching techniques in the presence of noise. Motivated by applications in learning theory, Awasthi et al. [5] considered recovery of a vector based on noisy 1-bit linear measurements. The resistance to noise demonstrated is analogous to what we show for the KOR sketch, but technically quite different since the linear mapping is computed over the reals before a sign operation is applied.

In a very recent paper [11], Choi et al. propose a framework for releasing differentially private estimates of various sketching problems in a distributed setting. This framework ensures that the estimates only have a multiplicative error factor. The technique relies on secure multi-party computation and the sketches submitted by each participant are not private and so cannot be released. Further, the results of Choi et al. do not immediately allow for estimating size or weight of the symmetric difference between two sets.

If the sketching matrix H itself is secret and randomly chosen from a distribution over matrices with entries in a finite field, very strong privacy guarantees on the sketch Hx can be obtained, while still allowing $\|x\|_0$ to be estimated from Hx with small error [7]. Blocki et al. [8] prove that the Johnson-Lindenstrauss transform is in fact differentially private, when keeping the sketch matrix secret. However, the condition that the sketch matrix is secret is a serious limitation for applications such as streaming and distributed cardinality estimation that require H to be stored or shared.

3 Preliminaries

We let $[n] = \{1, 2, \dots, n\}$ and let $U = [u]$ be the universe that the datasets are taken from.

For a set $A \subseteq U$, we let x_A denote the characteristic vector for A , defined as

$$(x_A)_j = \begin{cases} 1, & j \in A \\ 0, & \text{otherwise} \end{cases}.$$

We write w_A (or w_{x_A}) for the weight vector for input set A such that

$$w_A = x_A \circ w$$

for fixed, public weights $w_j \in (0, 1]$, and \circ denotes the Hadamard product.

For vector $x = (x_1, \dots, x_u)$ we define $\|x\|_p = \left(\sum_{j=1}^u x_j^p\right)^{1/p}$ as the p -norm of x . For $p = 0$, we define $\|x\|_0 = \sum_{j=1}^u \mathbf{1}[x_j \neq 0]$, often called the zero-"norm". F_0 denotes the 0th frequency moment and represents the number of distinct elements in a stream (or a set). Frequency moments are well-known from the streaming literature, see for example [4].

Our sketch Hx_A is comprised of $\log(u)$ "levels", $H_i x_A$ for $i = 0, \dots, \log(u) - 1$. We refer to Section 4.1 for a description of these levels. Let n denote the size of the binary vector representation of $H_i x_A$ for each i . Hence, the size of the noisy KOR sketch $Hx_A + \varphi$ is $\tau = n \log u$. Note that n is fixed and depends on the privacy parameter ϵ and the accuracy parameter β .

Finally, we assume that sets and vectors are stored in a sparse representation, such that we can list the non-zero entries in the input vector x in time $O(\|x\|_0)$.

3.1 Hashing-based subsampling

The sketch matrix H is defined by several hash functions. For simplicity, we assume access to an oracle representing random hash functions, namely, that we can sample a fully random hash function, and it can be evaluated in constant time. We do not store the hash function as part of our sketch, so the space for our sketch does not include space required for storing the hash function. We believe it is possible to replace these hash functions with concrete, efficient hash functions that can be stored in small space while preserving the asymptotic bounds on accuracy, but in order to focus on privacy aspects, we have not pursued this direction. Importantly, the differential privacy of our method holds for any choice of hash function and does not depend on the random oracle assumption.

To ensure that adding two sketches gives a sketch for the symmetric difference, it is necessary that both players sample the same elements for each H_i . To ensure coordinated sampling, we use a hash function, so the same elements from U are sampled by both players. We use the following (standard) subsampling technique: let \mathcal{S} be the family of all fully random hash functions from U into $[0, 1]$. Let $s \sim \mathcal{S}$ uniformly at random. We sample an element j from the input set at level $i = 0, \dots, \log(u) - 1$ if and only if $s(j) \in (w_j/2^{i+1}, w_j/2^i]$. We refer the reader to the survey of Woodruff [48] for more details on subsampling.

3.2 Differential Privacy

Differential privacy is a statistical property of the behavior of a mechanism [16]. The guarantee is that an adversary who observes the output of a differentially private mechanism will only obtain negligible information about the presence or absence of a particular item in the input data. Intuitively, a differentially private mechanism is almost insensitive to

the presence or absence of a single element, in the sense that the probability of observing a specific result should be almost the same for any two neighboring sets.

In Definition 3, we define differential privacy formally in terms of databases. In our application, the databases are sets, and thus *neighboring* means that one set is a subset of the other, and their sizes differ by 1.

► **Definition 3** (Differential Privacy [16]). *For $\varepsilon \geq 0$, a randomized mechanism \mathcal{M} is said to be ε -differentially private (or purely differentially private) if for any two neighboring databases, S and T – i.e., databases differing in a single entry – and for all $W \subseteq \text{Range}(\mathcal{M})$ it holds that*

$$\Pr [\mathcal{M}(S) \in W] \leq e^\varepsilon \cdot \Pr [\mathcal{M}(T) \in W].$$

For $\varepsilon \geq 0$ and $\delta \in [0, 1]$, a randomized mechanism \mathcal{M} is said to be (ε, δ) -differentially private (or approximately differentially private) if for any two neighboring databases, S and T , and for all $W \subseteq \text{Range}(\mathcal{M})$ it holds that

$$\Pr [\mathcal{M}(S) \in W] \leq e^\varepsilon \cdot \Pr [\mathcal{M}(T) \in W] + \delta.$$

We show in section 5.1 that our protocol obtains ε -differential privacy.

Our protocol for estimating the weight of the symmetric difference works in the *local* model of differential privacy, where each player adds noise to their own sketch. It uses the general technique of achieving privacy by adding noise according to *sensitivity* of a function [16]. We note that our sketch would *also* work in a model where vectors supplied by the users are combined using a black-box multi-party *secure aggregation* [22, 34]. In this setting, only the sketch for the symmetric difference would be released, and thus, only this sketch would need to be differentially private, meaning that less noise is required.

We can use the Laplace mechanism [16] to get differentially private estimates of the weights of the input sets. These estimates can be used together with an estimate for the weight of the symmetric difference to compute estimates for the union and the intersection of the two input sets with error that is of the same magnitude as the error for estimating the symmetric difference. For more details about differential privacy, we refer the reader to, for example, [18].

4 Techniques

4.1 Sketch Description

In this section, we describe the noisy KOR sketch in detail. The description is self-contained, but we refer the interested reader to [13] for more background on (linear) sketches. As mentioned, our sketch combines the techniques from [30] with hashing-based subsampling to achieve a sketch that is robust against adding noise, as long as we know how much noise was added.

We first give the intuition behind the $n \times u$ -matrices H_i , that our sketch H is comprised of: Suppose that we have a rough estimate \hat{E} of $\|w\|_1$, accurate within a constant factor. Then we can obtain a more precise estimate by sampling (using a hash function) a fraction n/\hat{E} of the elements, for some parameter n , and computing the sketch from [30] of size n for the sampled elements. This gives an approximation of the number of sampled elements, which in turn gives an approximation of $\|w\|_1$ with small relative error. Since we do not know $\|w\|_1$ within a constant factor – especially in the setting where we are interested in

the size of the symmetric difference – we use hashing-based subsampling to sample each element j from the input set with probability $w_j/2^{i+1}$ for $i = 0, \dots, \log(u) - 1$. Thus for each i , we sample elements corresponding to approximately a $1/2^{i+1}$ fraction of the weight and compute the sketch from [30] of size n for the sampled elements. For one of these i we are guaranteed to sample approximately a fraction $n/\|w\|_1$ of the input weight assuming that $\|w\|_1 > n$. For this i , we can obtain a precise estimate of $\|w\|_1$ from the sketch.

We now define H_i formally. We first describe the sketch from [30] as a linear sketch over $\text{GF}(2)$. Let \mathcal{F} be the family of all hash functions from universe U into $[n]$, and pick $h \sim \mathcal{F}$ uniformly at random. The hash function h uniquely defines an $n \times u$ -matrix K , where

$$K_{k,j} = \begin{cases} 1, & \text{if } h(j) = k \\ 0, & \text{otherwise} . \end{cases}$$

We combine this with the following sampling technique:

Let \mathcal{S} be the family of all hash functions from U to $[0, 1]$. Sample $s \sim \mathcal{S}$ uniformly at random. The hash function s defines a $u \times u$ -diagonal matrix S_i for each $i = 0, \dots, \log(u) - 1$, defined by

$$(S_i)_{j,j} = \begin{cases} 1, & \text{if } s(j) \in (w_j/2^{i+1}, w_j/2^i] \\ 0, & \text{otherwise} . \end{cases}$$

The matrix-vector product $S_i x$ represents subsample of input vector x , where we sample each element with probability $w_j/2^{i+1}$.

We are finally ready to define H_i as $H_i = K S_i$, which is an $n \times u$ -matrix over $\text{GF}(2)$. By definition:

$$(H_i)_{k,j} = \begin{cases} 1, & (h(j) = k) \wedge (s(j) \in (w_j/2^{i+1}, w_j/2^i]) \\ 0, & \text{otherwise} . \end{cases}$$

The KOR sketch can be represented as an $n \log(u) \times u$ -matrix H , formed by stacking $H_1, \dots, H_{\log(u)}$.

Let \mathcal{N}_ε be a distribution over vectors from $\{0, 1\}^{n \log(u)}$, where each entry is 1 independently with probability p . We show in Section 5.2 that it suffices to set $p = 1/(2 + \varepsilon)$. Sample the noise (or *perturbation*) vector $\varphi \sim \mathcal{N}_\varepsilon$ independently and uniformly at random. The *noisy* KOR sketch of x is then computed (over $\text{GF}(2)$) as:

$$Hx + \varphi.$$

4.2 Estimation

Next, we describe how to compute a weight estimate from a sketch $Hx + \varphi$. Let w be the weight vector associated with x . Let φ_i be the restriction of φ to the entries that are added to $H_i x$ when adding φ to Hx . To compute an estimate for $\|w\|_1$, for each $i = 0, \dots, \log(u) - 1$ count the number of 1s in $H_i x + \varphi_i$, $Z_i = \|H_i x + \varphi_i\|_0$ and compute the interval:

$$I_i = \begin{cases} [0, u] & \text{if } Z_i \geq (1 - \gamma)n/2 \\ \left[2^i n \ln \left(\frac{\frac{1}{2/\varepsilon+1}}{1 - \frac{2Z_i}{(1+\gamma)n}} \right), 2^i n \ln \left(\frac{\frac{1}{2/\varepsilon+1}}{1 - \frac{2Z_i}{(1-\gamma)n}} \right) \right] & \text{otherwise.} \end{cases} \quad (1)$$

where $\gamma < \frac{\beta-1/n}{7e^3(2/\varepsilon+1)}$. Compute the intersection $I = \bigcap_{i=0}^{\log(u)-1} I_i$ and check if the maximum value in I is within a factor $(1 + \eta)$ of the minimum value in I for

$$\eta = \frac{6\gamma \left(e^3 \left(\frac{2}{\varepsilon} + 1 \right) - 1 \right)}{1 + \gamma - 2\gamma \left(e^3 \left(\frac{2}{\varepsilon} + 1 \right) \right)} .$$

If that is the case, every element in I is a good estimate for $\|w\|_1$ (having relative error at most $(1 + \beta)$) with high probability. Otherwise, $\|w\|_1$ is small with high probability, and we let the estimate for $\|w\|_1$ be 0. We analyze the accuracy of this estimator in Section 5.

4.3 Application to symmetric difference

In this section, we describe a differentially private protocol to compute an estimate for the weight of the symmetric difference between sets held by two parties. First, we show that the sum of two noisy KOR sketches, $Hx_A + \varphi$ and $Hx_B + \psi$, is a noisy KOR sketch for the symmetric difference, $H(x_{A\Delta B}) + (\varphi + \psi)$, which has the same properties as $Hx_A + \varphi$ and $Hx_B + \psi$, but for $\varepsilon' < \varepsilon$ as more noise is added.

► **Lemma 4.** *Adding two noisy KOR sketches with perturbation vectors $\varphi \sim \mathcal{N}_\varepsilon$ and $\psi \sim \mathcal{N}_\varepsilon$, respectively, will yield a noisy KOR sketch for the symmetric difference of the input sets with noise $\varphi + \psi \sim \mathcal{N}_{\varepsilon'}$ for $\varepsilon' = \varepsilon^2/(2 + 2\varepsilon)$.*

Proof. Let x_A and x_B be the input vectors from each of the two players. Let H be as defined in Section 4.1, and define φ, ψ as the noise vectors for the noisy KOR sketches for x_A and x_B , respectively. We have (over $\text{GF}(2)$) that

$$\begin{aligned} (Hx_A + \varphi) + (Hx_B + \psi) &= (Hx_A + Hx_B) + (\varphi + \psi) \\ &= H(x_A + x_B) + (\varphi + \psi). \end{aligned}$$

This is exactly the noisy KOR sketch for the symmetric difference with perturbation $\varphi + \psi$. Note that we observe a 1 in an entry of $\varphi + \psi$ with probability $p' = p(1-p) + (1-p)p = 2p(1-p)$. We show in Section 5.2 that we can let $p = \frac{1}{2+\varepsilon}$. Observe that

$$p' = \frac{1}{2+\varepsilon'} = \frac{2}{2+\varepsilon} \left(1 - \frac{1}{2+\varepsilon}\right)$$

which implies that $\varepsilon' = \varepsilon^2/(2 + 2\varepsilon)$. ◀

By Lemma 4 we can treat a sketch for the symmetric difference exactly like a sketch for input vector x although with a different privacy parameter ε' . Hence, Theorem 1 gives us Corollary 2, restated here for convenience:

► **Corollary 2.** *For accuracy parameter $\beta > 0$, consider an ε -differentially private noisy KOR sketch for a set A and an ε -differentially private noisy KOR sketch for a set B , based on the same linear sketch $H \sim \mathcal{H}$, sampled independently of A and B . We can compute an approximation $\hat{\Delta}$ of the weight of the symmetric difference, such that with probability $1 - 1/u$:*

$$\| \|w_{A\Delta B}\|_1 - \hat{\Delta} \| < \beta \|w_{A\Delta B}\|_1 + \text{poly}(1/\varepsilon, 1/\beta, \log u) .$$

Note that the additive error in Corollary 2 still depends polynomially on ε even for privacy parameter ε' , which is explained by the fact that $\varepsilon' = \varepsilon^2/(2 + 2\varepsilon)$.

Finally, we assumed that $\|w_A\|_1$ and $\|w_B\|_1$ were released with Laplacian noise, which gives an expected additive error of $O(1/\varepsilon)$ for each of $\|w_A\|_1$ and $\|w_B\|_1$ [16]. We can use the following equations to get estimates for the union, intersection and difference:

$$\begin{aligned} \|w_{A\cup B}\|_1 &= \frac{\|w_A\|_1 + \|w_B\|_1 + \|w_{A\Delta B}\|_1}{2}, \\ \|w_{A\cap B}\|_1 &= \frac{\|w_A\|_1 + \|w_B\|_1 - \|w_{A\Delta B}\|_1}{2} \\ \|w_{A\setminus B}\|_1 &= \frac{\|w_A\|_1 + \|w_{A\Delta B}\|_1 - \|w_B\|_1}{2} . \end{aligned}$$

That is, the error is bounded by half the error of the estimate of the symmetric difference size plus $O(1/\varepsilon)$.

5 Proof of Theorem 1

In this section we give a proof of Theorem 1, restated here for convenience:

► **Theorem 1** (Noisy KOR sketch). *Let $w \in (0, 1]^u$ be given. For every choice of $0 < \beta < 1$ and $\varepsilon = O(1)$ there exists a distribution \mathcal{H} over $GF(2)$ -linear sketches mapping a vector $x \in \{0, 1\}^u$ to $\{0, 1\}^\tau$, where $\tau = O(\log^2(u)\varepsilon^{-2}\beta^{-2})$, and a distribution \mathcal{N}_ε over noise vectors such that:*

1. *For $H \sim \mathcal{H}$ and $\varphi \sim \mathcal{N}_\varepsilon$, given $Hx + \varphi$ we can compute, in time $O(\tau)$, an estimate \hat{w} of $\|x \circ w\|_1$ that with probability $1 - 1/u$ satisfies $|\hat{w} - \|x \circ w\|_1| < \beta \|x \circ w\|_1 + O(\log(u)\varepsilon^{-2}\beta^{-2})$.*
2. *For every H in the support of \mathcal{H} , $Hx + \varphi$ is ε -differentially private over the choice of $\varphi \sim \mathcal{N}_\varepsilon$, and can be computed in time $O(\|x\|_0 \log(u) + \tau)$, including time for sampling φ .*

5.1 Noise level and Differential Privacy Guarantees

We first show that the noisy KOR sketch $Hx + \varphi$ satisfies ε -differential privacy, which proves part 2 of Theorem 1. Intuitively, removal/insertion of a single element can change only a single entry in the sketch, as the element is inserted into only a single level.

► **Lemma 5.** *If $p \in \left(\frac{1}{e^\varepsilon + 1}, \frac{1}{2}\right)$ then $Hx + \varphi$ is ε -differentially private.*

Proof Sketch. The proof follows from the privacy of the Randomized Response Technique [46]. To make this work self-contained, we included a full proof in Appendix A.1. ◀

5.2 Bounding accuracy

In this section, along with Section 5.3, we prove the first part of Theorem 1. Let an input vector x be given and define w to be the corresponding weight vector. We will mainly consider each $H_i x$ isolated, so let φ_i be the n -dimensional (binary) randomness vector as described in the proof of Lemma 5. First, we state two useful lemmas.

► **Lemma 6.** *For each $i = 0, \dots, \log(u) - 1$ let $L_i = \|H_i x\|_0$ and $Z_i = \|H_i x + \varphi_i\|_0$. Then:*

$$\mathbb{E}_{\substack{h \sim \mathcal{F}, \\ s \sim \mathcal{S}}} [L_i] = \frac{n}{2} \left(1 - \prod_{j \in A} \left(1 - \frac{w_j}{2^i n} \right) \right) \quad (2)$$

$$\mathbb{E}_{\substack{h \sim \mathcal{F}, \\ s \sim \mathcal{S}, \\ \varphi_i \sim \mathcal{N}_p}} [Z_i] = \frac{n}{2} \left(1 - (1 - 2p) \prod_{j \in A} \left(1 - \frac{w_j}{2^i n} \right) \right) \quad (3)$$

Proof. We refer the reader to Appendix A.2 for the proof. ◀

► **Lemma 7.** *For $i = 0, \dots, \log(u) - 1$ let $Z_i = \|H_i x + \varphi_i\|_0$. For any $0 < \gamma < 1$, we have with probability at least $1 - 6 \log(u) e^{-\frac{\gamma^2 p^3 n}{6^2 \cdot 3}}$ that for all $i = 0, \dots, \log(u) - 1$ simultaneously:*

$$(1 - \gamma) \mathbb{E}_{\substack{h \sim \mathcal{F}, \\ s \sim \mathcal{S}, \\ \varphi_i \sim \mathcal{N}_p}} [Z_i] < Z_i < (1 + \gamma) \mathbb{E}_{\substack{h \sim \mathcal{F}, \\ s \sim \mathcal{S}, \\ \varphi_i \sim \mathcal{N}_p}} [Z_i].$$

Proof. We refer the reader to Appendix A.3 for the proof. \blacktriangleleft

First, we consider the case when $1 < n < \|w\|_1$. In Lemma 8 we state that in this case, with high probability we get an error of at most a factor $(1 + \beta)$ for a well-chosen γ , where γ is a function of the privacy parameter ε , the accuracy parameter β and the size of the universe, u . For convenience, define

$$I_i(p) = \begin{cases} [0, u] & \text{if } Z_i \geq (1 - \gamma)n/2 \\ \left[2^i n \ln \left(\frac{1-2p}{1 - \frac{2Z_i}{(1+\gamma)n}} \right), 2^i n \ln \left(\frac{1-2p}{1 - \frac{2Z_i}{(1-\gamma)n}} \right) \right] & \text{otherwise} \end{cases} \quad (4)$$

and $\hat{w} := 2^i n \ln \left(1 / \prod_{j \in A} \left(1 - \frac{w_j}{2^i n} \right) \right)$. We prove our result in two steps:

1. If $\hat{w} \in I_i(p)$ for all $i = 0, \dots, \log(u) - 1$, then there is some i such that any value from (4) estimates \hat{w} up to a factor $(1 + \eta)$, where η is a function of γ and ε .
2. $\|w\|_1 \leq \hat{w} \leq (1 + \frac{1}{2^i n}) \|w\|_1$ for each i . Specifically, $\|w\|_1 \leq \hat{w} \leq (1 + \frac{1}{n}) \|w\|_1$ for all i . Hence, we choose γ independent of i such that $(1 + \eta) (1 + \frac{1}{n}) \leq (1 + \beta)$ for at least one of the intervals $I_i(p)$. We pick γ to work for the i where $\|w\|_1 / (2^i n) \in [1, 2)$ as this corresponds to having an input of size between n and $2n$ (we obtain this input size by the sampling from x in H_i). If $\|w\|_1 \geq n$, there is such an i , and we can identify it by checking that the endpoints of the interval are sufficiently close together, as described in Section 4.2. We consider the case when $\|w\|_1 < n$ in Section 5.3 where we show that in this case, the error is bounded by an additive factor of $O(n)$.

► **Lemma 8.** *Assume $\|w\|_1 > n > 1$, and $\beta > \frac{1}{n}$. With probability at least $1 - 6 \log(u) e^{-\frac{\gamma^2 p^3 n}{108}}$ there exists an $i \in \{0, \dots, \log(u) - 1\}$ such that any element from $I_i(p)$ is a $(1 + \beta)$ -approximation to $\|w\|_1$ for*

$$\gamma < \frac{(\beta - \frac{1}{n})(1 - 2p)}{7e^3}.$$

Specifically, i where $\frac{\|w\|_1}{2^i n} \in [1, 2)$, gives these guarantees.

Proof Sketch. We give an informal sketch of the proof and refer the reader to Appendix A.4 for the formal proof. We first remark that for γ as described, Lemma 7 implies that if $\|w\|_1 / (2^i n) \leq 2$, then $Z_i < (1 - \gamma)n/2$ with high probability. Hence, it suffices to consider the intervals from (4) of the form $I_i(p) = \left[2^i n \ln \left(\frac{1-2p}{1 - \frac{2Z_i}{(1+\gamma)n}} \right), 2^i n \ln \left(\frac{1-2p}{1 - \frac{2Z_i}{(1-\gamma)n}} \right) \right]$. Define

$$\hat{w} := 2^i n \ln \left(\frac{1}{\prod_{j \in A} \left(1 - \frac{w_j}{2^i n} \right)} \right).$$

From Lemma 6, we have

$$\prod_{j \in A} \left(1 - \frac{w_j}{2^i n} \right) = \frac{1 - \frac{2\mathbb{E}[Z_i]}{n}}{1 - 2p}.$$

Assume that the bounds in Lemma 7 are satisfied. We remove this assumption shortly. By the bounds in Lemma 7, $\hat{w} \in I_i(p)$ for all i . We show that \hat{w} is contained in an interval, which is slightly bigger than $I_i(p)$ whenever $\|w\|_1 / (2^i n) \in [1, 2)$ and show that the endpoints of this interval are within a factor $(1 + \eta)$ of each other, where η is a function of γ . Clearly, then $I_i(p)$ is also sufficiently small for this i . Denote this interval $I_i^*(p)$. Any element from

$I_i^*(p)$ is a $(1+\eta)$ -approximation to \hat{w} . Removing the assumption that the bounds in Lemma 7 hold, we simply get a small error probability and conclude that with probability at least $1 - 6 \log(u) e^{-\gamma^2 p^3 n / 108}$ we have $\hat{w} \in I_i(p)$ for all i , and thus any value from $I_i^*(p)$ is a $(1+\eta)$ estimation to \hat{w} with high probability. Observing that $\|w\|_1 < \hat{w} \leq (1 + \frac{1}{n}) \|w\|_1$ for any i , we choose γ in terms of β such that $(1+\eta)(1 + \frac{1}{n}) < (1+\beta)$. Then any value from $I_i^*(p)$ is a $(1+\beta)$ -approximation for $\|w\|_1$. We formally choose γ in Appendix A.4. We remark that the assumption $\|w\|_1 / (2^i n) \in [1, 2)$ allows us to choose γ independent of i , such that we can compute $I_i(p)$ for all i with a single value of γ . \blacktriangleleft

Observing that $\frac{1}{2+\varepsilon} > \frac{1}{e^\varepsilon+1}$ for $\varepsilon > 0$, we let $p = 1/(2+\varepsilon)$ and observe that for $I_i := I_i(1/(2+\varepsilon))$ with the choice of γ described in Lemma 8, we get the interval I_i in (1).

5.3 Putting things together

In this section we consider the accuracy in the remaining case where $\|w\|_1 \leq n$. We also analyze the running time. Combining with Section 5.1 this completes the proof of Theorem 1.

Note that if $\varepsilon > 1$, we can start our protocol by dividing ε by a suitable constant, c such that $\varepsilon' = \varepsilon/c < 1$. Changing ε by a constant will change our bounds by a constant factor as well. Hence, we can without loss of generality assume $\varepsilon < 1$. We can also, without loss of generality, assume $u > 10$ – this will at most increase the failure probability and space by a constant factor.

We first show a sufficient upper bound on the sketch size $\tau = n \log u$. Observe that $p > 1/4$ and let $c_\gamma = 7e^3$ be a constant. Then we want $e^{-\frac{\gamma^2 p^3 n}{108}} < 1/u^2$ as this ensures a failure probability of at most $6 \log(u)/u^2 < 1/u$. Noting that

$$(1-2p)^2 = \left(1 - \frac{2}{2+\varepsilon}\right)^2 = \left(\frac{1}{2/\varepsilon+1}\right)^2 = \frac{1}{4/\varepsilon^2 + 4/\varepsilon + 1} > \frac{\varepsilon^2}{20},$$

we have

$$\begin{aligned} e^{-\frac{\gamma^2 p^3 n}{108}} &< e^{-\frac{\left(\frac{(\beta-\frac{1}{n})^{(1-2p)}}{7e^3}\right)^2 n/4^3}{108}} = e^{-\frac{(\beta-\frac{1}{n})^2 \left(\frac{1}{(2/\varepsilon+1)^2}\right)^n}{4^3 \cdot c_\gamma^2 \cdot 108}} \\ &< e^{-\frac{(\beta-\frac{1}{n})^2 \varepsilon^2 n}{20 \cdot 4^3 \cdot c_\gamma^2 \cdot 108}} < 1/u^2 \end{aligned}$$

when letting $n = O(\log(u)\beta^{-2}\varepsilon^{-2})$. Hence, the size of the sketch is

$$\tau = \log(u) \cdot n = O\left(\frac{\log^2(u)}{\varepsilon^2 \beta^2}\right).$$

Note that this n satisfies the requirement $\beta > 1/n$ from Lemma 8.

We argue about the error: Note that if $\|w\|_1 \geq n$, then if one of the intervals I_i is sufficiently small and $\hat{w} \in I_i$ for all $i = 0, \dots, \log(u) - 1$, then $\hat{w} \in I = \bigcap_{i=0}^{\log(u)-1} I_i$ and I is also sufficiently small to give the wanted estimate. So by Lemma 8, we can check if the endpoints of I are within a factor at most $(1+\eta)$ of each other, and if so, with probability $1 - 1/u$ any value from I is within a factor $(1+\beta)$ of $\|w\|_1$. If I is too big, then none of the intervals I_i was sufficiently small implying that our assumption that $\|w\|_1 / (2^i n) \in [1, 2)$ does not hold for any i . Hence, with probability $1 - 1/u$ we have $\|w\|_1 < n$. We refer to the formal proof in Appendix A.4 for the details. Our protocol sets the estimate of $\|w\|_1$ to 0 leading to an additive error of $O(n)$ when I was too big. This means that we get an additive

error of at most $n = O(\log(u)\beta^{-2}\varepsilon^{-2})$, as required.

Finally, we comment on the running times: For the first part of Theorem 1, we note that in order to compute the estimate, we need to count the number of ones in $H_i x + \varphi_i$ for each $i = 0, \dots, \log(u) - 1$, compute the intervals I_i and their intersection and check if it is sufficiently small. Counting the number of ones in all $H_i x + \varphi_i$ is the bottleneck and requires time $O(\tau)$. For the second part of Theorem 1, note that we can initialize the randomness vector φ in time $O(\tau)$ and we can hash vector x in time $O(\|x\|_0 \log(u))$ assuming that we can iterate over x in time $O(\|x\|_0)$.

Combining with Lemma 8 and Lemma 5, we have completed the proof of Theorem 1.

6 Distributed Streaming Implementation

In a streaming setting want a sketch which can be updated and two sketches can be merged to give a sketch for the union of the input streams, while we cannot guarantee that there are no duplicates in the input stream. In this case, our sketch does not immediately apply, as items with an even number of occurrences would "cancel out". Such items would therefore never be represented in the sketch, as the sketch is over $\text{GF}(2)$. This issue can easily be fixed: the idea is to add another layer of sampling, such that we sample each *occurrence* of a data item with probability $1/2$. Hence, we treat identical items independently on each occurrence and so ensures that an entry in the sketch is 1 with probability $1/2$, regardless of the number of copies of identical items and collisions with other items. We refer to this as the *pre-sampled* sketch. The intuition is that the number of copies of an item inserted in the pre-sampled sketch is even or odd with probability $1/2$. By Chernoff bounds the fraction of elements that are sampled an odd number of times is very close to $1/2$ with high probability. Thus it is natural to consider the estimator that is two times the estimator described in Section 4.2.

To understand this in more detail we argue that merging two (non-private) pre-sampled sketches over $\text{GF}(2)$ gives a sketch for the *union* of the two input sets. Suppose $z \in A \cup B$, $h(z) = k$ and that z is sampled at level i . We argue that $\Pr[(H_i x_{A \cup B})_k = 1] = 1/2$. Note that

$$(H_i x_{A \cup B})_k = 1 \quad \Leftrightarrow \quad (H_i x_A)_k \neq (H_i x_B)_k.$$

Further, we have that if $z \in A$, then $\Pr[(H_i x_A)_k = 1] = 1/2$ regardless of the number of other elements hashing to k at level i . If no elements from A hash to entry k at level i , then $\Pr[(H_i x_A)_k = 1] = 0$. We have

$$\begin{aligned} \Pr[(H_i x_{A \cup B})_k = 1] &= \Pr[(H_i x_A)_k = 1] \Pr[(H_i x_B)_k = 0] \\ &\quad + \Pr[(H_i x_A)_k = 0] \Pr[(H_i x_B)_k = 1], \end{aligned}$$

which is $1/2$ whenever $z \in A \cup B$.

Acknowledgement. We thank Shuang Song and Abhradeep Guha Thakurta for feedback on a previous version of this manuscript.

References

- 1 Mohammad Alaggan, Sébastien Gambs, and Anne-Marie Kermarrec. BLIP: non-interactive differentially-private similarity computation on bloom filters. In *Stabilization, Safety, and Security of Distributed Systems - 14th International Symposium, SSS*, pages 202–216, 2012. doi:10.1007/978-3-642-33536-5_20.
- 2 Mohammad Alaggan, Sébastien Gambs, Stan Matwin, and Mohammed Tuhin. Sanitization of call detail records via differentially-private bloom filters. In *Data and Applications Security and Privacy XXIX - 29th Annual IFIP WG 11.3 Working Conference, DBSec 2015*, pages 223–230, 2015. doi:10.1007/978-3-319-20810-7_15.
- 3 Noga Alon, Phillip B Gibbons, Yossi Matias, and Mario Szegedy. Tracking join and self-join sizes in limited storage. *Journal of Computer and System Sciences*, 64(3):719–747, 2002.
- 4 Noga Alon, Yossi Matias, and Mario Szegedy. The space complexity of approximating the frequency moments. In *Symposium on the Theory of Computing*, pages 20–29, 1996. doi:10.1145/237814.237823.
- 5 Pranjal Awasthi, Maria-Florina Balcan, Nika Haghtalab, and Hongyang Zhang. Learning and 1-bit compressed sensing under asymmetric noise. In *Conference on Learning Theory*, pages 152–192, 2016.
- 6 Ziv Bar-Yossef, TS Jayram, Ravi Kumar, D Sivakumar, and Luca Trevisan. Counting distinct elements in a data stream. In *International Workshop on Randomization and Approximation Techniques in Computer Science*, pages 1–10, 2002.
- 7 Valerio Bioglio, Tiziano Bianchi, and Enrico Magli. Secure compressed sensing over finite fields. In *International Workshop on Information Forensics and Security (WIFS)*, pages 191–196, 2014.
- 8 Jeremiah Blocki, Avrim Blum, Anupam Datta, and Or Sheffet. The johnson-lindenstrauss transform itself preserves differential privacy. In *Symposium on Foundations of Computer Science, FOCS*, pages 410–419, 2012. doi:10.1109/FOCS.2012.67.
- 9 Andrei Z. Broder and Michael Mitzenmacher. Survey: Network applications of bloom filters: A survey. *Internet Mathematics*, 1(4):485–509, 2003. doi:10.1080/15427951.2004.10129096.
- 10 Clément Canonne, Gautam Kamath, and Thomas Steinke. The discrete gaussian for differential privacy. *arXiv preprint arXiv:2004.00010*, 2020.
- 11 Seung Geol Choi, Dana Dachman-Soled, Mukul Kulkarni, and Arkady Yerukhimovich. Differentially-private multi-party sketching for large-scale statistics. *IACR Cryptol. ePrint Arch.*, 2020:29, 2020. URL: <https://eprint.iacr.org/2020/029>.
- 12 Reuven Cohen, Liran Katzir, and Aviv Yehezkel. A unified scheme for generalizing cardinality estimators to sum aggregation. *Information Processing Letters*, 115(2):336–342, 2015.
- 13 Graham Cormode, Minos N. Garofalakis, Peter J. Haas, and Chris Jermaine. Synopses for massive data: Samples, histograms, wavelets, sketches. *Foundations and Trends in Databases*, 4(1-3):1–294, 2012. doi:10.1561/19000000004.
- 14 Damien Desfontaines, Andreas Lochbihler, and David A. Basin. Cardinality estimators do not preserve privacy. *PoPETs*, 2019(2):26–46, 2019. doi:10.2478/popets-2019-0018.
- 15 Devdatt P Dubhashi and Alessandro Panconesi. *Concentration of measure for the analysis of randomized algorithms*. Cambridge University Press, 2009.
- 16 Cynthia Dwork, Frank McSherry, Kobbi Nissim, and Adam D. Smith. Calibrating noise to sensitivity in private data analysis. In *3rd Theory of Cryptography Conference, TCC*, pages 265–284, 2006. doi:10.1007/11681878_14.
- 17 Cynthia Dwork, Moni Naor, Toniann Pitassi, Guy N Rothblum, and Sergey Yekhanin. Pan-private streaming algorithms. In *ICS*, pages 66–80, 2010.
- 18 Cynthia Dwork and Aaron Roth. The algorithmic foundations of differential privacy. *Foundations and Trends in Theoretical Computer Science*, 9(3-4):211–407, 2014. doi:10.1561/04000000042.

- 19 Úlfar Erlingsson, Vasyl Pihur, and Aleksandra Korolova. RAPPOR: randomized aggregatable privacy-preserving ordinal response. In *Proceedings of the 2014 Conference on Computer and Communications Security*, pages 1054–1067, 2014. doi:10.1145/2660267.2660348.
- 20 Philippe Flajolet, Éric Fusy, Olivier Gandouet, and Frédéric Meunier. HyperLogLog: the analysis of a near-optimal cardinality estimation algorithm. In *AofA: Analysis of Algorithms*, pages 137–156, 2007.
- 21 Philippe Flajolet and G. Nigel Martin. Probabilistic counting algorithms for data base applications. *J. Comput. Syst. Sci.*, 31(2):182–209, 1985. doi:10.1016/0022-0000(85)90041-8.
- 22 Slawomir Goryczka, Li Xiong, and Vaidy S. Sunderam. Secure multiparty aggregation with differential privacy: a comparative study. In *Joint 2013 EDBT/ICDT Conferences, EDBT/ICDT '13*, pages 155–163, 2013. doi:10.1145/2457317.2457343.
- 23 Peter J Haas, Jeffrey F Naughton, S Seshadri, and Lynne Stokes. Sampling-based estimation of the number of distinct values of an attribute. In *VLDB*, volume 95, pages 311–322, 1995.
- 24 Moritz Hardt and Kunal Talwar. On the geometry of differential privacy. In *Symposium on Theory of Computing, STOC*, pages 705–714, 2010. doi:10.1145/1806689.1806786.
- 25 T. S. Jayram and David P. Woodruff. Optimal bounds for Johnson-Lindenstrauss transforms and streaming problems with subconstant error. *Transactions on Algorithms*, 9(3):26:1–26:17, 2013. doi:10.1145/2483699.2483706.
- 26 Daniel M Kane, Jelani Nelson, and David P Woodruff. An optimal algorithm for the distinct elements problem. In *Proceedings of the 29th ACM symposium on Principles of database systems (PODS)*, pages 41–52, 2010.
- 27 Krishnamurthy Kenthapadi, Aleksandra Korolova, Ilya Mironov, and Nina Mishra. Privacy via the Johnson-Lindenstrauss transform. *J. Priv. Confidentiality*, 5(1), 2013. doi:10.29012/jpc.v5i1.625.
- 28 Daniel Kifer, Shai Ben-David, and Johannes Gehrke. Detecting change in data streams. In *VLDB*, volume 4, pages 180–191. Toronto, Canada, 2004.
- 29 Daniel Kifer and Ashwin Machanavajjhala. No free lunch in data privacy. In *Proceedings of ACM International Conference on Management of data (SIGMOD)*, pages 193–204, 2011.
- 30 Eyal Kushilevitz, Rafail Ostrovsky, and Yuval Rabani. Efficient search for approximate nearest neighbor in high dimensional spaces. In *Symposium on the Theory of Computing*, pages 614–623, 1998. doi:10.1145/276698.276877.
- 31 Andrew McGregor, Ilya Mironov, Toniann Pitassi, Omer Reingold, Kunal Talwar, and Salil Vadhan. The limits of two-party differential privacy. In *51st Annual Symposium on Foundations of Computer Science*, pages 81–90, 2010.
- 32 Frank McSherry and Kunal Talwar. Mechanism design via differential privacy. In *FOCS*, volume 7, pages 94–103, 2007.
- 33 Frank D McSherry. Privacy integrated queries: an extensible platform for privacy-preserving data analysis. In *Proceedings of ACM International Conference on Management of data (SIGMOD)*, pages 19–30, 2009.
- 34 Luca Melis, George Danezis, and Emiliano De Cristofaro. Efficient private statistics with succinct sketches. In *23rd Annual Network and Distributed System Security Symposium, NDSS*, 2016. doi:10.14722/ndss.2016.23175.
- 35 Darakhshan Mir, S Muthukrishnan, Aleksandar Nikolov, and Rebecca N Wright. Pan-private algorithms: When memory does not help. *arXiv preprint arXiv:1009.1544*, 2010.
- 36 Darakhshan Mir, Shan Muthukrishnan, Aleksandar Nikolov, and Rebecca N Wright. Pan-private algorithms via statistics on sketches. In *Proceedings of the 30th Symposium on Principles of Database Systems (PODS)*, pages 37–48, 2011.
- 37 Ilya Mironov. On significance of the least significant bits for differential privacy. In Ting Yu, George Danezis, and Virgil D. Gligor, editors, *Conference on Computer and Communications Security, CCS*, pages 650–661, 2012. doi:10.1145/2382196.2382264.

- 38 Ilya Mironov, Omkant Pandey, Omer Reingold, and Salil P. Vadhan. Computational differential privacy. In Shai Halevi, editor, *Advances in Cryptology - CRYPTO*, volume 5677 of *Lecture Notes in Computer Science*, pages 126–142, 2009. doi:10.1007/978-3-642-03356-8_8.
- 39 Michael Mitzenmacher, Rasmus Pagh, and Ninh Pham. Efficient estimation for high similarities using odd sketches. In *Proceedings of 23rd international conference on World Wide Web (WWW)*, pages 109–118, 2014.
- 40 Aleksandar Nikolov. Personal communication. 2020.
- 41 Hagen Sparka, Florian Tschorsch, and Björn Scheuermann. P2KMV: A privacy-preserving counting sketch for efficient and accurate set intersection cardinality estimations. *IACR Cryptology ePrint Archive*, 2018:234, 2018.
- 42 Rade Stanojevic, Mohamed Nabeel, and Ting Yu. Distributed cardinality estimation of set operations with differential privacy. In *IEEE Symposium on Privacy-Aware Computing, PAC*, pages 37–48, 2017. doi:10.1109/PAC.2017.43.
- 43 Florian Tschorsch and Björn Scheuermann. An algorithm for privacy-preserving distributed user statistics. *Computer Networks*, 57(14):2775–2787, 2013. doi:10.1016/j.comnet.2013.05.011.
- 44 Salil P. Vadhan. The complexity of differential privacy. In *Tutorials on the Foundations of Cryptography*, pages 347–450. Springer, 2017. doi:10.1007/978-3-319-57048-8_7.
- 45 Saskia Nuñez von Voigt and Florian Tschorsch. Rtxfm: Probabilistic counting for differentially private statistics. In *Workshop on Trust and Privacy Aspects of Smart Information Environments (TPSIE)*, 2019.
- 46 Stanley L. Warner. Randomized response: A survey technique for eliminating evasive answer bias. *Journal of the American Statistical Association*, 60(309):63–69, 1965. URL: <http://www.jstor.org/stable/2283137>.
- 47 Royce J Wilson, Celia Yuxin Zhang, William Lam, Damien Desfontaines, Daniel Simmons-Marengo, and Bryant Gipson. Differentially private SQL with bounded user contribution. *Proceedings on Privacy Enhancing Technologies*, 2020(2):230–250, 2020.
- 48 David P. Woodruff. Data streams and applications in computer science. *Bulletin of the EATCS*, 114, 2014. URL: <http://eatcs.org/beatcs/index.php/beatcs/article/view/304>.

A Omitted proofs

A.1 Differential Privacy Guarantees

► **Lemma 5.** *If $p \in \left(\frac{1}{e^\epsilon + 1}, \frac{1}{2}\right)$ then $Hx + \varphi$ is ϵ -differentially private.*

Proof. Let A and B be two neighboring input sets with corresponding characteristic vectors, x_A and x_B , where neighboring means that one set is a subset of the other and the sizes differ by 1. By symmetry of differential privacy, we can without loss of generality assume that A is the smaller set. Suppose that $B \setminus \{z\} = A$. The element z can only affect $H_i x$ for i where z is sampled. If z is never sampled, then $Hx_A = Hx_B$ and privacy is trivial. So assume $i \in \{0, \dots, \log(u) - 1\}$ such that $s(z) \in (w_z/2^{i+1}, w_z/2^i]$. We limit our attention to $H_i x_A + \varphi_i$, where we can think of φ_i as the restriction of the $n \log(u)$ -dimensional random vector $\varphi \sim \mathcal{N}_\epsilon$ to the entries that would be added to $H_i x_A$ when adding φ to Hx_A . We show that $H_i x_A + \varphi_i$ is ϵ -differentially private. This implies that the entire sketch, $Hx_A + \varphi$, is ϵ -differentially private.

Inserting z into the sketch implies that $H_i x_A$ and $H_i x_B$ will differ in exactly one entry, i.e., $\|H_i x_A + H_i x_B\|_0 = 1$. Fix a noisy sketch, S_i . There exist unique vectors φ_i and ψ_i , such that $S_i = H_i x_A + \varphi_i = H_i x_B + \psi_i$. Note that $\|\varphi_i - \psi_i\|_0 = 1$. Let $\|\varphi_i\|_0 = r$. Then $\|\psi_i\|_0 = r'$ for $r' \in \{r + 1, r - 1\}$. Conditioned on $\|\varphi_i\|_0 = r$ and $\|\psi_i\|_0 = r'$, the probabilities of randomly drawing exactly these randomness vectors are, respectively:

$$(1 - p)^{n-r} p^r \quad \text{and} \quad (1 - p)^{n-r'} p^{r'}.$$

Let $\varepsilon = O(1)$ be given. By Section 3.2 it is enough to show that for any fixed output $S_i = H_i x_A + \varphi_i = H_i x_B + \psi_i$, we have

$$e^{-\varepsilon} \leq \frac{\Pr[\text{observe } S_i \text{ from } A]}{\Pr[\text{observe } S_i \text{ from } B]} = \frac{\Pr[\text{observe } H_i x_A + \varphi_i \text{ from } A]}{\Pr[\text{observe } H_i x_B + \psi_i \text{ from } B]} \leq e^\varepsilon.$$

where the probability is over the randomness in φ_i and ψ_i . The sketches for A and B are computed using the same H_i , so the choice of H_i has no impact.

Hence, to obtain differential privacy it suffices that for every possible value of r and $r' \in \{r+1, r-1\}$

$$e^{-\varepsilon} \leq \frac{(1-p)^{n-r} p^r}{(1-p)^{n-r'} p^{r'}} = \frac{1}{(1-p)^{r-r'} p^{r'-r}} \leq e^\varepsilon,$$

which is satisfied for $1/2 > p \geq 1/(e^\varepsilon + 1)$, since $p < 1/2$ by assumption. \blacktriangleleft

A.2 Expectations

► **Lemma 6.** For each $i = 0, \dots, \log(u) - 1$ let $L_i = \|H_i x\|_0$ and $Z_i = \|H_i x + \varphi_i\|_0$. Then:

$$\mathbb{E}_{\substack{h \sim \mathcal{F}, \\ s \sim \mathcal{S}}} [L_i] = \frac{n}{2} \left(1 - \prod_{j \in A} \left(1 - \frac{w_j}{2^i n} \right) \right) \quad (2)$$

$$\mathbb{E}_{\substack{h \sim \mathcal{F}, \\ s \sim \mathcal{S}, \\ \varphi_i \sim \mathcal{N}_p}} [Z_i] = \frac{n}{2} \left(1 - (1-2p) \prod_{j \in A} \left(1 - \frac{w_j}{2^i n} \right) \right) \quad (3)$$

Proof. Let A be the input set with corresponding weight vector w . Let $v_i \in \mathbb{Z}_{\geq 0}^n$ be a vector such that for each $k \in [n]$

$$(v_i)_k = \sum_{j \in A} \mathbf{1} \left[\frac{s(j)}{w_j} \in (1/2^{i+1}, 3/2^{i+1}] \right] \cdot \mathbf{1}[h(j) = k].$$

That is, each entry $(v_i)_k$ is the number of candidates for entry k in the sketch at level i , i.e., the number of items j that hash to k and satisfy $\frac{s(j)}{w_j} \in (1/2^{i+1}, 3/2^{i+1}]$. Since $s(j)$ is uniform, we have for such a candidate

$$\Pr_{s \sim \mathcal{S}} \left[s(j) \in (w_j/2^{i+1}, 2w_j/2^{i+1}] \mid s(j) \in (w_j/2^{i+1}, 3w_j/2^{i+1}] \right] = \frac{1}{2}.$$

If there is at least one candidate for entry k then, by the Principle of Deferred Decisions, the probability that we sample an odd number of these is $1/2$ and so for $i = 0, \dots, \log(u) - 1$

$$\begin{aligned} \Pr_{\substack{h \sim \mathcal{F}, \\ s \sim \mathcal{S}}} [(H_i x_A)_k = 1 \mid (v_i)_k \neq 0] &= \frac{1}{2}, \\ \Pr_{\substack{h \sim \mathcal{F}, \\ s \sim \mathcal{S}}} [(H_i x_A)_k = 1 \mid (v_i)_k = 0] &= 0. \end{aligned}$$

As

$$\Pr_{s \sim \mathcal{S}} \left[\frac{s(j)}{w_j} \in (1/2^{i+1}, 3/2^{i+1}] \right] = \Pr_{s \sim \mathcal{S}} [s(j) \in (w_j/2^{i+1}, 3w_j/2^{i+1})] = \frac{w_j}{2^i},$$

we have

$$\Pr_{\substack{h \sim \mathcal{F} \\ s \sim \mathcal{S}}} [(v_i)_k \neq 0] = 1 - \prod_{j \in A} \left(1 - \frac{w_j}{2^i n}\right).$$

We conclude that

$$\Pr_{\substack{h \sim \mathcal{F} \\ s \sim \mathcal{S}}} [(H_i x_A)_k = 1] = \frac{1 - \prod_{j \in A} \left(1 - \frac{w_j}{2^i n}\right)}{2}.$$

and letting $L_i = \sum_{k=1}^n (H_i x_A)_k$, we get

$$\mathbb{E}_{\substack{h \sim \mathcal{F} \\ s \sim \mathcal{S}}} [L_i] = \frac{n}{2} \left(1 - \prod_{j \in A} \left(1 - \frac{w_j}{2^i n}\right)\right)$$

We similarly compute an expression for $\mathbb{E}_{h \sim \mathcal{F}, s \sim \mathcal{S}, \varphi_i \sim \mathcal{N}_p} [Z_i]$. Let φ_i be the restriction of a randomness vector $\varphi \sim \mathcal{N}_\varepsilon$ to the entries that are added to $H_i x_A$ when adding φ to $H x_A$. We see that

$$\begin{aligned} & \Pr_{\substack{h \sim \mathcal{F} \\ s \sim \mathcal{S} \\ \varphi_i \sim \mathcal{N}_p}} [(H_i x_A + \varphi_i)_k = 1] \\ &= \Pr_{\substack{h \sim \mathcal{F} \\ s \sim \mathcal{S} \\ \varphi_i \sim \mathcal{N}_p}} [(H_i x_A + \varphi_i)_k = 1 \mid (H_i x_A)_k = 1] \cdot \Pr_{\substack{h \sim \mathcal{F} \\ s \sim \mathcal{S}}} [(H_i x_A)_k = 1] \\ &+ \Pr_{\substack{h \sim \mathcal{F} \\ s \sim \mathcal{S} \\ \varphi_i \sim \mathcal{N}_p}} [(H_i x_A + \varphi_i)_k = 1 \mid (H_i x_A)_k = 0] \cdot \Pr_{\substack{h \sim \mathcal{F} \\ s \sim \mathcal{S}}} [(H_i x_A)_k = 0] \\ &= (1-p) \cdot \Pr_{\substack{h \sim \mathcal{F} \\ s \sim \mathcal{S}}} [(H_i x_A)_k = 1] + p \cdot \Pr_{\substack{h \sim \mathcal{F} \\ s \sim \mathcal{S}}} [(H_i x_A)_k = 0] \\ &= (1-p) \cdot \frac{1}{2} \left(1 - \prod_{j \in A} \left(1 - \frac{w_j}{2^i n}\right)\right) + p \cdot \left(1 - \frac{1 - \prod_{j \in A} \left(1 - \frac{w_j}{2^i n}\right)}{2}\right) \\ &= \frac{1}{2} - \left(\frac{1}{2} - p\right) \prod_{j \in A} \left(1 - \frac{w_j}{2^i n}\right) \end{aligned}$$

showing that

$$\mathbb{E}_{\substack{h \sim \mathcal{F} \\ s \sim \mathcal{S} \\ \varphi_i \sim \mathcal{N}_p}} [Z_i] = \frac{n}{2} \left(1 - (1-2p) \prod_{j \in A} \left(1 - \frac{w_j}{2^i n}\right)\right).$$

◀

A.3 Concentration bounds

► **Lemma 7.** For $i = 0, \dots, \log(u) - 1$ let $Z_i = \|H_i x + \varphi_i\|_0$. For any $0 < \gamma < 1$, we have with probability at least $1 - 6 \log(u) e^{-\frac{\gamma^2 p^3 n}{6^2 \cdot 3}}$ that for all $i = 0, \dots, \log(u) - 1$ simultaneously:

$$(1 - \gamma) \mathbb{E}_{\substack{h \sim \mathcal{F} \\ s \sim \mathcal{S} \\ \varphi_i \sim \mathcal{N}_p}} [Z_i] < Z_i < (1 + \gamma) \mathbb{E}_{\substack{h \sim \mathcal{F} \\ s \sim \mathcal{S} \\ \varphi_i \sim \mathcal{N}_p}} [Z_i].$$

Before proving Lemma 7, we mention the following lemma:

► **Lemma 9.** *Let $L_i = \|H_i x_A\|_0$. For any $0 < \gamma' < 1$, we have with probability at least $1 - 4 \log(u) e^{-2\gamma'^2 n}$*

$$\mathbb{E}_{\substack{h \sim \mathcal{F}, \\ s \sim \mathcal{S}}} [L_i] - 2\gamma'n \leq L_i \leq \mathbb{E}_{\substack{h \sim \mathcal{F}, \\ s \sim \mathcal{S}}} [L_i] + 2\gamma'n$$

for all $i = 0, \dots, \log(u) - 1$ simultaneously.

Proof. Let A be the input set and w the corresponding weight vector. Let $v_i \in \mathbb{Z}_{\geq 0}^n$ be a vector such that for each $k \in [n]$

$$(v_i)_k = \sum_{j \in A} \mathbf{1} \left[\frac{s(j)}{w_j} \in (1/2^{i+1}, 3/2^{i+1}] \right] \cdot \mathbf{1} [h(j) = k]$$

so $(v_i)_k$ is the number of candidates for entry k in the sketch at level i . Let $V_i = \|v_i\|_0 = \sum_{k=1}^n \mathbf{1}[(v_i)_k \neq 0]$. V_i is a sum of negatively associated random variables (for definition and argument see Section 4.1 in [15]), so by Theorem 4.3 in [15], we can use the Hoeffding bound to see that with probability at least $1 - 2e^{-2n\gamma'^2}$ we have for any $i = 0, \dots, \log(u) - 1$

$$\mathbb{E}[V_i] - \gamma'n \leq V_i \leq \mathbb{E}[V_i] + \gamma'n. \quad (5)$$

Let $L_i = \|H_i x_A\|_0 = \sum_{k=1}^n (H_i x_A)_k$ denote the number of ones in the linear sketch. For fixed V_i , L_i is a sum of independent random variables with (by the principle of deferred decisions)

$$\Pr \left[(H_i x_A)_k = 1 \mid (v_i)_k \neq 0 \right] = \frac{1}{2}, \quad \Pr \left[(H_i x_A)_k = 1 \mid (v_i)_k = 0 \right] = 0.$$

So for any fixed $V_i = t$

$$\mathbb{E} \left[L_i \mid V_i = t \right] = \frac{t}{2}. \quad (6)$$

Furthermore, as L_i is a sum of independent random variables for a fixed choice of V_i , we can use the Hoeffding bound: with probability at least $1 - 2e^{-2n\gamma'^2}$

$$\mathbb{E} \left[L_i \mid V_i = t \right] - \gamma'n \leq L_i|_{V_i=t} \leq \mathbb{E} \left[L_i \mid V_i = t \right] + \gamma'n,$$

where $L_i|_{V_i=t}$ means the value of L_i when we assume that $V_i = t$. Combining this with (5) and (6) a union bound gives with probability at least $1 - 4e^{-2n\gamma'^2}$

$$\frac{\mathbb{E}[V_i] - \gamma'n}{2} - \gamma'n \leq L_i \leq \frac{\mathbb{E}[V_i] + \gamma'n}{2} + \gamma'n. \quad (7)$$

Simultaneously, (5) and (6) gives

$$\frac{\mathbb{E}[V_i] - \gamma'n}{2} \leq \mathbb{E}[L_i] \leq \frac{\mathbb{E}[V_i] + \gamma'n}{2}, \quad (8)$$

which implies

$$2\mathbb{E}[L_i] - \gamma'n \leq \mathbb{E}[V_i] \leq 2\mathbb{E}[L_i] + \gamma'n. \quad (9)$$

Note that in the union bound from (7), we already assumed that (5) was satisfied, so (9) is trivially satisfied under the union bound without changing the probability guarantees. Hence, inserting (9) into (7), we have

$$\frac{2\mathbb{E}[L_i] - 2\gamma'n}{2} - \gamma'n \leq L_i \leq \frac{2\mathbb{E}[L_i] + 2\gamma'n}{2} + \gamma'n. \quad (10)$$

which finally shows that with probability at least $1 - 4e^{-2n\gamma'^2}$ we have

$$\mathbb{E}[L_i] - 2\gamma'n \leq L_i \leq \mathbb{E}[L_i] + 2\gamma'n.$$

A union bound over the $\log(u)$ values of i concludes the proof. \blacktriangleleft

We are now ready to prove Lemma 7.

Proof of Lemma 7. Fix i . Let $L_i = \|H_i x_A\|_0$ and $Z_i = \|H_i x_A + \varphi_i\|_0$. We let $Z_i|_{L_i=t}$ be the number of ones in $H_i x_A + \varphi_i$, assuming that $L_i = t$. For any fixed value $t \in \{0, \dots, n\}$ of L_i , we have

$$\mathbb{E}_{\varphi_i \sim \mathcal{N}_p} [Z_i|_{L_i=t}] = (1-p) \cdot t + p(n-t) = np + (1-2p)t. \quad (11)$$

By Lemma 9, with probability at least $1 - 4\log(u)e^{-2\gamma'^2 n}$ we have for all $i = 0, \dots, \log(u) - 1$

$$\mathbb{E}_{\substack{h \sim \mathcal{F}, \\ s \sim \mathcal{S}, \\ \varphi_i \sim \mathcal{N}_p}} [Z_i] \geq np + (1-2p) \left(\mathbb{E}_{\substack{h \sim \mathcal{F}, \\ s \sim \mathcal{S}}} [L_i] - 2\gamma'n \right) \quad (12)$$

$$\mathbb{E}_{\substack{h \sim \mathcal{F}, \\ s \sim \mathcal{S}, \\ \varphi_i \sim \mathcal{N}_p}} [Z_i] \leq np + (1-2p) \left(\mathbb{E}_{\substack{h \sim \mathcal{F}, \\ s \sim \mathcal{S}}} [L_i] + 2\gamma'n \right) \quad (13)$$

Furthermore, for any fixed H_i , let $Z_i|_{H_i}$ denote the number of ones in $H_i x_A + \varphi_i$, conditioned on this choice of H_i . We note that fixing H_i is equivalent to fixing L_i as L_i is uniquely determined by H_i and the input. $Z_i|_{H_i}$ is a sum of independent random variables, where the randomness comes from the perturbation. So for any $0 < \gamma^* < 1$, a Chernoff bound gives

$$\Pr_{\varphi_i \sim \mathcal{N}_p} \left[Z_i|_{H_i} > (1 + \gamma^*) \mathbb{E} [Z_i|_{H_i}] \vee Z_i|_{H_i} < (1 - \gamma^*) \mathbb{E} [Z_i|_{H_i}] \right] \quad (14)$$

$$\leq 2e^{-\gamma^{*2} \mathbb{E} [Z_i|_{H_i}] / 3} \quad (15)$$

where $\mathbb{E} [Z_i|_{H_i}]$ is over $\varphi_i \sim \mathcal{N}_p$. By (11), $\mathbb{E}_{\varphi_i \sim \mathcal{N}_p} [Z_i|_{H_i}] \geq np$ for any choice of H_i , so $2e^{-\gamma^{*2} pn/3}$ is an upper bound on (15). Moreover, (15) holds for all $i = 0, \dots, \log(u) - 1$ simultaneously with probability at most $2\log(u)e^{-\gamma^{*2} pn/3}$. We conclude that

$$\Pr_{\varphi_i \sim \mathcal{N}_p} \left[\forall i : (1 - \gamma^*) \mathbb{E} [Z_i|_{H_i}] < Z_i|_{H_i} < (1 + \gamma^*) \mathbb{E} [Z_i|_{H_i}] \right] \quad (16)$$

$$\geq 1 - 2\log(u)e^{-\gamma^{*2} pn/3} \quad (17)$$

Combining (12), (13) and (17) and letting $\gamma' = \gamma^*$, we have by a union bound that for all levels i simultaneously, where the expectation is over $h \sim \mathcal{F}$ and $s \sim \mathcal{S}$

$$\begin{aligned} Z_i &\geq (1 - \gamma') (np + (1 - 2p) (\mathbb{E}[L_i] - 2\gamma'n)) \\ Z_i &\leq (1 + \gamma') (np + (1 - 2p) (\mathbb{E}[L_i] + 2\gamma'n)), \end{aligned}$$

with probability at least

$$1 - \left(4 \log(u) e^{-2n\gamma'^2} + 2 \log(u) e^{-\gamma'^2 pn/3}\right) \geq 1 - 6 \log(u) e^{-\gamma'^2 pn/3}.$$

By Lemma 6, this is equivalent to

$$Z_i \geq (1 - \gamma') (E[Z_i] - 2(1 - 2p)\gamma' n) \quad (18)$$

$$Z_i \leq (1 + \gamma') (E[Z_i] + 2(1 - 2p)\gamma' n). \quad (19)$$

where the expectation is over $h \sim \mathcal{F}$, $s \sim \mathcal{S}$ and $\varphi_i \sim \mathcal{N}_p$. We pick a suitable γ' :

$$\begin{aligned} \gamma' = \frac{\gamma p}{6} &\Rightarrow 2(1 - 2p)\gamma' n = (1 - 2p) \frac{\gamma p}{3} n \\ &\Rightarrow 2(1 - 2p)\gamma' n \leq \frac{\gamma(1 - 2p)}{3} E[Z_i]. \end{aligned}$$

Hence, let $\gamma' = \frac{\gamma p}{6}$. Inserting into (18) and (19) we have

$$Z_i \geq \left(1 - \frac{\gamma p}{6}\right) \left(E[Z_i] - \frac{\gamma(1 - 2p)}{3} E[Z_i]\right)$$

$$Z_i \leq \left(1 + \frac{\gamma p}{6}\right) \left(E[Z_i] + \frac{\gamma(1 - 2p)}{3} E[Z_i]\right)$$

where $E[Z_i]$ is over $h \sim \mathcal{F}$, $s \sim \mathcal{S}$ and $\varphi_i \sim \mathcal{N}_p$.

We conclude that with this choice of γ , with probability at least $1 - 6 \log(u) e^{-\frac{\gamma^2 p^3 n}{6^2 \cdot 3}}$

$$(1 - \gamma) \underset{\substack{h \sim \mathcal{F}, \\ s \sim \mathcal{S}, \\ \varphi_i \sim \mathcal{N}_p}}{E} [Z_i] \leq Z_i \leq (1 + \gamma) \underset{\substack{h \sim \mathcal{F}, \\ s \sim \mathcal{S}, \\ \varphi_i \sim \mathcal{N}_p}}{E} [Z_i].$$

◀

A.4 Size of interval for input size

Before proving Lemma 8, we give a technical lemma:

► **Lemma 10.** For any $0 < \gamma < \frac{1}{\frac{2e^3}{1-2p} - 1}$ any value

$$\hat{m} \in \left[2^i n \ln \left(\frac{1 - 2p}{1 - \frac{2Z_i}{(1+\gamma)n}} \right), 2^i n \ln \left(\frac{1 - 2p}{1 - \frac{2Z_i}{(1-\gamma)n}} \right) \right] \quad (20)$$

satisfies

$$\hat{m} \geq (1 - \eta) 2^i n \ln \left(\frac{1}{\prod_{j \in A} \left(1 - \frac{w_j}{2^i n}\right)} \right)$$

$$\hat{m} \leq (1 + \eta) 2^i n \ln \left(\frac{1}{\prod_{j \in A} \left(1 - \frac{w_j}{2^i n}\right)} \right)$$

for

$$\eta = \frac{6\gamma \left(\frac{e^3}{1-2p} - 1 \right)}{(1 - \gamma) - 2\gamma \left(\frac{e^3}{1-2p} - 1 \right)}$$

with probability at least $1 - 6 \log(u) e^{-\gamma^2 p^3 n/108}$ for the i where $\frac{\|w\|_1}{2^i n} \in [1, 2]$.

Proof. By Lemma 6

$$\prod_{j \in A} \left(1 - \frac{w_j}{2^i n}\right) = \frac{1 - \frac{2\mathbb{E}[Z_i]}{n}}{1 - 2p}$$

and so by Lemma 7, with probability at least $1 - 6 \log(u) e^{-\gamma^2 p^3 n / 108}$ we have for any $0 < \gamma < 1$ that for all $i = 0, \dots, \log(u) - 1$ simultaneously.

$$\frac{1 - \frac{2Z_i}{(1-\gamma)n}}{1 - 2p} < \prod_{j \in A} \left(1 - \frac{w_j}{2^i n}\right) < \frac{1 - \frac{2Z_i}{(1+\gamma)n}}{1 - 2p}. \quad (21)$$

For convenience, we consider the slightly bigger interval – note that if (21) is satisfied, then so is this interval:

$$\frac{1 - \frac{2(1+\gamma)\mathbb{E}[Z_i]}{(1-\gamma)n}}{1 - 2p} < \prod_{j \in A} \left(1 - \frac{w_j}{2^i n}\right) < \frac{1 - \frac{2(1-\gamma)\mathbb{E}[Z_i]}{(1+\gamma)n}}{1 - 2p},$$

where the left-hand side can be reordered as

$$\left(1 - \frac{2\gamma}{1-\gamma} \left(\frac{1}{(1-2p) \prod_{j \in A} \left(1 - \frac{w_j}{2^i n}\right)} - 1\right)\right) \prod_{j \in A} \left(1 - \frac{w_j}{2^i n}\right) \quad (22)$$

and the right-hand side as

$$\left(1 + \frac{2\gamma}{1+\gamma} \left(\frac{1}{(1-2p) \prod_{j \in A} \left(1 - \frac{w_j}{2^i n}\right)} - 1\right)\right) \prod_{j \in A} \left(1 - \frac{w_j}{2^i n}\right). \quad (23)$$

We will bound this interval further using the following claim:

▷ **Claim 11.** Define

$$\beta^* := \frac{2\gamma}{1-\gamma} \left(\frac{e^{2+\frac{1}{2^i-1n}}}{1-2p} - 1\right).$$

Whenever $\frac{\|w\|_1}{2^i n} < 2$, the interval defined by (22) and (23) is contained in

$$\left[(1 - \beta^*) \prod_{j \in A} \left(1 - \frac{w_j}{2^i n}\right), (1 + \beta^*) \prod_{j \in A} \left(1 - \frac{w_j}{2^i n}\right) \right]$$

Proof of Claim. As $\frac{2\gamma}{1+\gamma} < \frac{2\gamma}{1-\gamma}$, we increase (23) to

$$\left(1 + \frac{2\gamma}{1-\gamma} \left(\frac{1}{(1-2p) \prod_{j \in A} \left(1 - \frac{w_j}{2^i n}\right)} - 1\right)\right) \prod_{j \in A} \left(1 - \frac{w_j}{2^i n}\right).$$

Observing that when $\frac{\|w\|_1}{2^i n} \leq 2$

$$\begin{aligned} \frac{2\gamma}{1-\gamma} \left(\frac{1}{(1-2p) \prod_{j \in A} \left(1 - \frac{w_j}{2^i n}\right)} - 1\right) &\leq \frac{2\gamma}{1-\gamma} \left(\frac{e^{\frac{\|w\|_1}{2^i n} + \frac{\|w\|_1}{(2^i n)^2}}}{1-2p} - 1\right) \\ &\leq \frac{2\gamma}{1-\gamma} \left(\frac{e^{2+\frac{1}{2^i-1n}}}{1-2p} - 1\right) =: \beta^* \end{aligned}$$

we have the result. ◀

Applying the claim, we consider the interval:

$$2^i n \ln \left(\frac{1}{\prod_{j \in A} \left(1 - \frac{w_j}{2^i n}\right)} \right) \geq 2^i n \ln \left(\frac{1}{(1 + \beta^*) \prod_{j \in A} \left(1 - \frac{w_j}{2^i n}\right)} \right) \quad (24)$$

$$2^i n \ln \left(\frac{1}{\prod_{j \in A} \left(1 - \frac{w_j}{2^i n}\right)} \right) \leq 2^i n \ln \left(\frac{1}{(1 - \beta^*) \prod_{j \in A} \left(1 - \frac{w_j}{2^i n}\right)} \right). \quad (25)$$

We remind the reader that by construction, this interval contains the target interval (20).

We consider the ratio between the end-points of the interval defined by (24) and (25).

Observe that

$$\begin{aligned} \frac{2^i n \ln \left(\frac{1}{(1 - \beta^*) \prod_{j \in A} \left(1 - \frac{w_j}{2^i n}\right)} \right)}{2^i n \ln \left(\frac{1}{(1 + \beta^*) \prod_{j \in A} \left(1 - \frac{w_j}{2^i n}\right)} \right)} &= \frac{\ln \left(\frac{1}{\prod_{j \in A} \left(1 - \frac{w_j}{2^i n}\right)} \right) - \ln(1 - \beta^*)}{\ln \left(\frac{1}{\prod_{j \in A} \left(1 - \frac{w_j}{2^i n}\right)} \right) - \ln(1 + \beta^*)} \\ &\leq \frac{\ln \left(\frac{1}{\prod_{j \in A} \left(1 - \frac{w_j}{2^i n}\right)} \right) + \frac{\beta^*}{1 - \beta^*}}{\ln \left(\frac{1}{\prod_{j \in A} \left(1 - \frac{w_j}{2^i n}\right)} \right) - \beta^*} \\ &= 1 + \frac{\beta^* \left(1 + \frac{1}{1 - \beta^*}\right)}{\ln \left(\frac{1}{\prod_{j \in A} \left(1 - \frac{w_j}{2^i n}\right)} \right) - \beta^*} \end{aligned}$$

where the inequality follows from

$$\frac{x}{1+x} \leq \ln(1+x) \leq x, \quad x > -1.$$

For $\beta^* < 1/2$, we have

$$\begin{aligned} \frac{\beta^* \left(1 + \frac{1}{1 - \beta^*}\right)}{\ln \left(\frac{1}{\prod_{j \in A} \left(1 - \frac{w_j}{2^i n}\right)} \right) - \beta^*} &< \frac{3\beta^*}{\ln \left(\frac{1}{\prod_{j \in A} \left(1 - \frac{w_j}{2^i n}\right)} \right) - \beta^*} \\ &< \frac{3\beta^*}{\frac{\|w\|_1}{2^i n} - \beta^*} \end{aligned}$$

Observe that as $\frac{\|w\|_1}{2^i n}$ increases, it gets easier to satisfy this inequality. But we remind ourselves of the Claim, where we required $\frac{\|w\|_1}{2^i n} < 2$. So the interval in (24) and (25) does not necessarily contain the target interval (20) for larger values of $\frac{\|w\|_1}{2^i n}$. Assume further that $\frac{\|w\|_1}{2^i n} \geq 1$. Then

$$\frac{3\beta^*}{\frac{\|w\|_1}{2^i n} - \beta^*} < \frac{3\beta^*}{1 - \beta^*}.$$

So, we conclude that with probability at least $1 - 6 \log(u) e^{-\gamma^2 p^3 n / 108}$, any value in the target interval (20) is within a factor $1 + \frac{3\beta^*}{1 - \beta^*}$ of $2^i n \ln \left(\left(\prod_{j \in A} \left(1 - \frac{w_j}{2^i n}\right) \right)^{-1} \right)$.

Inserting the value of β^* , we obtain an estimate within a factor of

$$1 + \frac{6\gamma \left(\frac{e^{2 + \frac{1}{2^i - 1n}}}{1 - 2p} - 1 \right)}{(1 - \gamma) - 2\gamma \left(\frac{e^{2 + \frac{1}{2^i - 1n}}}{1 - 2p} - 1 \right)} < 1 + \frac{6\gamma \left(\frac{e^3}{1 - 2p} - 1 \right)}{(1 - \gamma) - 2\gamma \left(\frac{e^3}{1 - 2p} - 1 \right)}.$$

Thus it suffices that

$$\gamma < \frac{1}{\frac{2e^3}{1-2p} - 1}.$$

◀

We are now ready to prove Lemma 8:

► **Lemma 8.** *Assume $\|w\|_1 > n > 1$, and $\beta > \frac{1}{n}$. With probability at least $1 - 6 \log(u) e^{-\frac{\gamma^2 p^3 n}{108}}$ there exists an $i \in \{0, \dots, \log(u) - 1\}$ such that any element from $I_i(p)$ is a $(1 + \beta)$ -approximation to $\|w\|_1$ for*

$$\gamma < \frac{(\beta - \frac{1}{n})(1 - 2p)}{7e^3}.$$

Specifically, i where $\frac{\|w\|_1}{2^i n} \in [1, 2)$, gives these guarantees.

Proof. We will choose γ in terms of the accuracy parameter β , such that with high probability any estimate from the interval

$$\left[2^i n \ln \left(\frac{1 - 2p}{1 - \frac{2Z_i}{(1+\gamma)n}} \right), 2^i n \ln \left(\frac{1 - 2p}{1 - \frac{2Z_i}{(1-\gamma)n}} \right) \right] \quad (26)$$

is within a factor $(1 + \beta)$ of $\|w\|_1$. We do this in a few steps: First, we show that any value from (26) is a good estimate of

$$2^i n \ln \left(\frac{1}{\prod_{j \in A} \left(1 - \frac{w_j}{2^i n} \right)} \right). \quad (27)$$

As $2^i n \ln \left(\frac{1}{e^{-\frac{\|w\|_1}{2^i n}}} \right) = \|w\|_1$ and

$$\frac{2^i n \ln \left(\frac{1}{\prod_{j \in A} \left(1 - \frac{w_j}{2^i n} \right)} \right)}{2^i n \ln \left(\frac{1}{e^{-\frac{\|w\|_1}{2^i n}}} \right)} \leq \frac{\ln \left(e^{\frac{\|w\|_1}{2^i n} + \frac{\|w\|_1}{(2^i n)^2}} \right)}{\ln \left(e^{\frac{\|w\|_1}{2^i n}} \right)} = 1 + \frac{1}{2^i n}$$

where we used the Taylor expansion of the exponential function, we have

$$\|w\|_1 \leq 2^i n \ln \left(\frac{1}{\prod_{j \in A} \left(1 - \frac{w_j}{2^i n} \right)} \right) \leq \left(1 + \frac{1}{2^i n} \right) \|w\|_1.$$

So a good estimate for (27) will allow for a good estimate of $\|w\|_1$. The technical lemma, Lemma 10, shows that as long as $\|w\|_1$ is sufficiently large, that is, there is an i such that $\frac{\|w\|_1}{2^i n} \in [1, 2)$, we get a suitable estimate for (27) with the interval (26) with high probability.

Hence, any value from (26) is within a factor $(1 + \beta)$ of $\|w\|_1$ for

$$\gamma < \frac{(\beta - 1/n)(1 - 2p)}{7e^3} < \frac{\beta - 1/n}{7 \left(\frac{e^3}{1-2p} - 1 \right)} < \frac{\beta - \frac{1}{2^i n}}{7 \left(\frac{e^3}{1-2p} - 1 \right)}$$

for $\beta > \frac{1}{n}$. We will choose n in terms of β such that this is always satisfied. Clearly, this value of γ is significantly smaller than the requirement from Lemma 10, which concludes the proof. ◀