
A Distributed Algorithm For Differentially Private Heavy Hitters

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 In this paper we present a simple *distributed* algorithm for learning the heavy
2 hitters in a stream of data in the local model of differential privacy. Our algorithm
3 is efficient in the sense that its running time and communication complexity is
4 polynomial in the input size. Furthermore, our algorithm achieves the optimal error
5 guarantee for the frequency estimation.

6 1 Introduction

7 One of the most widely performed learning tasks on telemetry data by software companies is to
8 understand *heavy hitters or frequent items*. Typical applications of such a data analysis include
9 understanding the most frequently visited websites (URLs), popular apps on an app-store, trending
10 news, etc. In order to meet users' privacy expectations and in view of tightening privacy regulations (European
11 GDPR law) the ability to collect telemetry data privately is paramount. In this paper we are
12 interested in algorithms that preserve user's privacy in the face of such data collection from users.

13 Over the past decade, Differential privacy, introduced in the seminal work of Dwork *et al* [12]
14 (DP), has emerged as the standard definition of privacy for data analysis and machine learning, with
15 wide acceptance both in academia and industry. In the context of telemetry collection one typically
16 considers algorithms that exhibit a stronger notion of differential privacy called *local differential*
17 *privacy* [14, 18, 16, 10, 5], also known as randomized response model [21], γ -amplification [15],
18 or FRAPP [1]. These are randomized algorithms that are invoked on *user's device* to turn user's
19 private value into a response that is communicated to data collector and have the property that the
20 likelihood of any specific algorithm's output varies little with the input, thus providing users with
21 plausible deniability. Formally, given a constant $\epsilon > 0$, a randomized algorithm $\mathcal{A} : \mathcal{V} \rightarrow \mathcal{Z}$ is
22 ϵ -locally differentially private (ϵ -LDP) if for any pair of values $v, v' \in \mathcal{V}$ and any subset of output
23 $S \subseteq \mathcal{Z}$, we have that $\Pr[\mathcal{A}(v) \in S] \leq e^\epsilon \cdot \Pr[\mathcal{A}(v') \in S]$.

24 In the heavy hitters problem we are given a set of n users; each user holds an item v_i , which we assume
25 without loss of generality is an element in $\{1, 2, \dots, d\}$. A histogram of the data lists for each item
26 $v \in [d]$ the fraction of population that holds item v . However, when d is large compared to n , as is the
27 case in our example of understanding most frequent URLs, an explicit representation of histogram is
28 neither possible nor necessary. Often it is sufficient to have an implicit representation of the histogram.
29 Two commonly used representations are: a) *Frequency oracle*, which is an algorithm that answers the
30 frequency of a queried item $v \in [d]$. The frequency of an item v is defined as $f(v) = \frac{1}{n} \cdot |\{i | v_i = v\}|$.
31 b) *Heavy hitters*, which returns a list of k most frequent items. An algorithm for heavy items returns
32 a list of k items $\hat{v}_1, \hat{v}_2, \dots, \hat{v}_k$ together with the estimated frequencies $\hat{f}(\hat{v}_1), \hat{f}(\hat{v}_2), \dots, \hat{f}(\hat{v}_k)$. The
33 frequencies of items not in the list are assumed to be zero. The error of a heavy hitter algorithm for
34 item v is defined to be $|\hat{f}(v) - f(v)|$. The error of heavy hitter algorithm is the maximum error over
35 all the items $v \in [d]$. That is, Error = $\max_{v \in [d]} |\hat{f}(v) - f(v)|$.

36 The heavy hitters problem in the context of local differential privacy was made popular in the
 37 influential paper of Erlingsson, Pihur, and Korolova [14]. They gave differentially private heuristics,
 38 without any rigorous error analysis, for the heavy hitters problem that formed the basis of its
 39 implementation in Chrome browser. Later, Apple used differentially private heavy hitters algorithms
 40 to understand most frequently used emojis and to discover new words in text data [20].

41 The theoretical study of the problem, however, dates back to [17, 19], who gave private algorithms
 42 for heavy hitters with error $O(\frac{1}{n^{1/6}})$. The algorithms of [17] only satisfied (ϵ, δ) -differential privacy
 43 for $\delta > 0$, which is a strictly weaker notion of differential privacy. In a breakthrough work Bassily
 44 and Smith [4] constructed *optimal* algorithms for differentially private heavy hitters problem. They
 45 showed that there exist efficient mechanisms for both frequency oracle and heavy hitters with an error
 46 $O(\sqrt{\log d/n})$.

47 Despite its theoretical optimality on the error, the algorithm in [4] is not very suitable for use in
 48 practice. The client side running time of the algorithm is $O(n^{1.5})$, which is prohibitively expensive
 49 in applications where the number of users is in the order of hundreds of millions. This limitation
 50 was addressed in a very recent paper [2], who gave two elegant algorithms, *TreeHist* and *Bitstogram*,
 51 which dramatically reduced the running time on the client side to $O(\max\{\log n, \log d\}^2)$. Both these
 52 algorithms rely on first solving the frequency oracle problem, and then using it as a subroutine to
 53 solve the heavy hitters problem. The Bitstogram algorithm achieves the optimal guarantees on the
 54 error, whereas the TreeHist algorithm is sub-optimal by a factor of $\sqrt{\log n}$ for the heavy hitters
 55 problem.

56 All the known algorithms for the heavy hitters [4, 2] problem rely on users having access to *shared*
 57 *random bits*. In the TreeHist algorithm every user needs access to $O(\log n)$ number of shared
 58 random hash functions from a universal family of hash functions. On the other hand, Bitstogram
 59 algorithm requires that users have access to a shared random $\{-1, 1\}$ -matrix of size $O(\sqrt{n} \times n)$. The
 60 most secure way (and sometimes the only way) to implement shared randomness is that the *server*
 61 *generates random bits and sends it to each client*. For Bitstogram algorithm, this would require
 62 transmitting $\Omega(\sqrt{n})$ number of bits for each round of data collection. (Note that a user does not
 63 need the entire matrix but only a single column.) This is too expensive as telemetry data is collected
 64 for many applications, and downloading $\Omega(\sqrt{n})$ number of bits when n is in order of hundreds of
 65 millions is just not feasible. On the other hand, TreeHist protocol needs access to $O(\log n)$ bits, and
 66 hence is more amenable in practice. Yet, both these schemes require that *all the users have same*
 67 *shared random bits*.

68 Implementing shared randomness, *however small it may be*, poses several engineering challenges
 69 for large companies with billions of users that have huge telemetry clients, which log thousands
 70 of telemetry events. To make matters worse, telemetry data is collected periodically (say every
 71 day or every 6 hours [9]), and most implementations also use *memoization* tables to reduce privacy
 72 leakage across different rounds [14, 9]. This means that the server needs to ensure that every client
 73 is *synchronized*, in the sense that they have the access to fresh random bits for every round of
 74 telemetry collection or when the memoization tables are flushed. This entails the server and clients
 75 maintaining a state for every single telemetry event and every round, and guaranteeing that every
 76 client is synchronized with the server. For many large scale telemetry systems this extra overhead can
 77 lead to several tricky implementation issues. Hence, based on our conversations with practitioners,
 78 the following question arose.

79 *"Is there a distributed algorithm for learning heavy hitters that eliminates the need for shared*
 80 *randomness across users?"*

81 This intriguing question is the motivation behind this paper. First we formalize what we mean by a
 82 distributed algorithm.

83 **Definition 1.** *An algorithm for heavy hitters problem is distributed if there is no shared random bits*
 84 *across users. Each user knows only the privacy parameter ϵ and size of the universe d .*

85 Note that the size of universe is often implicitly known in most applications. For example, in the case
 86 of URLs it is determined by the maximum size a URL can take. Furthermore, the privacy parameter
 87 is often configured in the client software. All known algorithms [4, 2] for heavy hitters do not satisfy
 88 the above definition.

89 **Our results.**

90 The main result of the paper is the following.

Theorem 1. *There exists an efficient and distributed algorithm, which we name as AoN, for the frequency oracle and the heavy hitters problems in the local model of differential privacy. AoN gives following guarantees: 1) AoN is ϵ -locally differentially private for any $\epsilon > 0$; 2) AoN achieves the optimal worst case error guarantee for the frequency oracle problem. The worst case error guarantee of the frequency oracle is at most*

$$\frac{e^{\frac{\epsilon}{2}}(1 + e^{\frac{\epsilon}{2}})}{e^{\frac{\epsilon}{2}} - 1} \cdot \sqrt{\frac{1}{2n} \log\left(\frac{2d}{\delta}\right)},$$

91 *which holds with probability $1 - \delta$; 3) The worst case error guarantee of the heavy hitters algo-*
 92 *rithm is at most $\frac{e^{\frac{\epsilon}{4}}(1+e^{\frac{\epsilon}{4}})}{e^{\frac{\epsilon}{4}} - 1} \cdot \sqrt{\frac{1}{2n} \log\left(\frac{2d}{\delta}\right) \log n}$, which holds with probability $1 - \delta$. 4) The client*
 93 *side running time of AoN is $O(\max\{\log n, \log d\}^2)$; 5) The communication complexity of AoN is*
 94 *$O(\max\{\log n, \log d\})$ bits.*

95 The error guarantees obtained by AoN matches that of TreeHist. We note that the number of bits
 96 transmitted from a client to the server is more in AoN compared to TreeHist ($O(\log n)$ bits compared
 97 $O(1)$). However, the number bits sent from the server to a client (or alternatively, the number of
 98 bits a client downloads from the server) is zero where as in TreeHist it $O(\log n)$. Moreover, AoN
 99 eliminates the need for the clients to be synchronized, which we believe is the novel aspect of our
 100 scheme compared to every existing LDP mechanism. We give a complete proof this theorem in the
 101 full version of the paper.

102 Given an optimal frequency oracle, we also show that heavy hitters problem can also be solved in a
 103 distributed manner by some modifications to the schemes in [2].

104 Apart from the properties mentioned in Theorem (1), our algorithm has several other desirable
 105 characteristics: 1) AoN is simple to specify and memoize compared to the existing mechanisms. 2)
 106 The constants in the error guarantee given by AoN are considerably smaller compared to the TreeHist
 107 algorithm. However, we would like to emphasize that the main technical contribution of the paper is
 108 a novel distributed and optimal algorithm for the frequency oracle problem with a theoretical analysis.
 109 To the best of our knowledge this is the first such known algorithm for the heavy hitters problem.

110 2 All-or-Nothing (AoN) Mechanism

111 Now we present AoN mechanism, which is a computation and communication efficient algorithm for
 112 computing heavy hitters. *An algorithm is computation and communication efficient if a) its running*
 113 *time is polynomial in n and $\log d$, b) the number bits transmitted by the algorithm is polynomial in n*
 114 *and $\log d$.* Observe that precisely because of this reason we can not simply use a histogram algorithm
 115 to solve the frequency estimation problem for heavy hitters problem. Such an approach would require
 116 m bits of communication. AoN mechanism consists of two parts: 1) A client side protocol that is
 117 invoked on each user device to report differentially private version of the item $v_i \in \{1, 2 \dots d\}$ that
 118 user i holds. 2) A frequency oracle on the server side that collects user reports and produces an
 119 estimated frequency of an item $v_i \in \{1, 2 \dots d\}$. Now we describe these two algorithms and prove
 120 their properties.

121 2.1 AoN Client Protocol

122 During the data collection, each user $i \in [n]$ samples uniformly at random a hash function $h_i \in H$
 123 from a universal family of hash functions H . (It can be shown that sampling a hash function from
 124 an universal family is equivalent to selecting two integers $a, b \in [p]$ uniformly at random.) Next,
 125 the user i evaluates $h_i(v_i)$. If $h_i(v_i) = 1$, then the user reports the hash function h_i to the server. If
 126 $h_i(v_i) \neq 1$, then user tosses a coin with bias $\frac{1}{e^\epsilon}$, where ϵ is the privacy parameter; if it lands heads,
 127 then the user sends h_i to the server. Otherwise, the user transmits an empty string ϕ . The formal
 128 description of the algorithm is given in Algorithm 1.

129 The below lemma follows immediately from the definition.

130 **Lemma 1.** *Each user $i \in [n]$ transmits $O(\max\{\log n, \log d\})$ bits to the server, and the running*
 131 *time of Algorithm 1 is $O(\max\{\log n, \log d\}^2)$.*

Algorithm 1 AoN Client Protocol

AoN Mechanism for Data Collection on Client

Input: an element $v_i \in [d]$ from user i , privacy parameter $\epsilon \geq 1$;

Output: a hash function $h_i \in H$ or ϕ , where each $h_i : [d] \rightarrow [B]$, and $B = \lceil e^{\epsilon/2} + 1 \rceil$;

- 1: Randomly draw a hash function h_i from a universal family $H = \{h : [d] \rightarrow [B]\}$;
 - 2: If $h_i(v_i) = 1$: output h_i ;
 - 3: Else: output h_i with probability $1/e^\epsilon$, and ϕ otherwise.
-

Algorithm 2 AoN Frequency Oracle

Processing Frequency Query on Server

Input: an element $v \in [d]$;

Output: an estimation $\hat{f}(v)$ of the frequency of v ($f(v) = |\{i \in [n] \mid v_i = v\}|/n$);

- 1: For each user $i \in [n]$:
- 2: Collect $A(v_i)$ using the client-side mechanism;
- 3: If $A(v_i) \neq \phi$:
- 4: If $h_i(v) = 1$: $\theta \leftarrow \theta + 1$;
- 5: Output the estimated frequency of v as

$$\hat{f}(v) = \left(\frac{\theta}{n} - \frac{1}{e^{\frac{\epsilon}{2}}(1 + e^{\frac{\epsilon}{2}})} \right) \cdot \frac{e^{\frac{\epsilon}{2}}(1 + e^{\frac{\epsilon}{2}})}{e^{\frac{\epsilon}{2}} - 1}$$

132 Next, we show that AoN is ϵ -differentially private. We defer the proof to a full version of the paper.

133 **Lemma 2.** *For any $\epsilon > 0$, Algorithm 1 is ϵ -differentially private.*

134 **2.2 AoN Frequency Oracle**

135 Our server side frequency oracle algorithm is quite simple. Suppose we want to estimate the frequency
136 of an item $v \in [d]$. The algorithm scans through hash functions received from the users, and counts
137 the number of hash functions that evaluate to 1 on the item v . Let θ be this quantity. Our algorithm
138 outputs the frequency of item v as

$$\hat{f}(v) = \left(\frac{\theta}{n} - \frac{1}{e^{\frac{\epsilon}{2}}(1 + e^{\frac{\epsilon}{2}})} \right) \cdot \frac{e^{\frac{\epsilon}{2}}(1 + e^{\frac{\epsilon}{2}})}{e^{\frac{\epsilon}{2}} - 1} \quad (1)$$

139 A formal description of our algorithm is given in Algorithm 2.

140 It remains to analyze the error guarantee of our algorithm. As a first step towards that we show that
141 the expected value of $\hat{f}(v)$ is equal to the true frequency of item v .

142 **Lemma 3.** *For every item $v \in [d]$, $\hat{f}(v)$ is an unbiased estimator of $f(v)$; that is $E[\hat{f}(v)] = f(v)$.*

143 We do a concentration analysis to prove the error guarantee.

144 **Lemma 4.** *For every item $v \in [d]$ and $\delta > 0$, with probability $1 - \delta$,*

$$|\hat{f}(v) - f(v)| \leq \frac{e^{\frac{\epsilon}{2}}(1 + e^{\frac{\epsilon}{2}})}{e^{\frac{\epsilon}{2}} - 1} \cdot \sqrt{\frac{1}{2n} \log \left(\frac{2d}{\delta} \right)}$$

145 The proof of the main theorem (1) now follows from Lemmas (1,2,3,4). The omitted proofs will be
146 made available in the full version of the paper.

147 **References**

- 148 [1] S. Agrawal and J. R. Haritsa. A framework for high-accuracy privacy-preserving mining. In
149 *ICDE*, pages 193–204, 2005.

- [2] R. Bassily, K. Nissim, U. Stemmer, and A. Thakurta. Practical locally private heavy hitters. In *NIPS*, 2017.
- [3] R. Bassily, K. Nissim, U. Stemmer, and A. G. Thakurta. Practical locally private heavy hitters. In *Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017, 4-9 December 2017, Long Beach, CA, USA*, pages 2285–2293, 2017.
- [4] R. Bassily and A. D. Smith. Local, private, efficient protocols for succinct histograms. In *STOC*, pages 127–135, 2015.
- [5] R. Bassily, A. D. Smith, and A. Thakurta. Private empirical risk minimization: Efficient algorithms and tight error bounds. In *FOCS*, pages 464–473, 2014.
- [6] S. Boucheron, G. Lugosi, and P. Massart. *Concentration inequalities: A nonasymptotic theory of independence*. Oxford university press, 2013.
- [7] J. L. Carter and M. N. Wegman. Universal classes of hash functions (extended abstract). In *Proceedings of the Ninth Annual ACM Symposium on Theory of Computing*, STOC ’77, pages 106–112, New York, NY, USA, 1977. ACM.
- [8] M. Dietzfelbinger, J. Gil, Y. Matias, and N. Pippenger. Polynomial hash functions are reliable (extended abstract). In *Proceedings of the 19th International Colloquium on Automata, Languages and Programming*, ICALP ’92, pages 235–246, London, UK, UK, 1992. Springer-Verlag.
- [9] B. Ding, J. Kulkarni, and S. Yekhanin. Collecting telemetry data privately. In *Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017, 4-9 December 2017, Long Beach, CA, USA*, pages 3574–3583, 2017.
- [10] J. C. Duchi, M. I. Jordan, and M. J. Wainwright. Local privacy and statistical minimax rates. In *FOCS*, pages 429–438, 2013.
- [11] J. C. Duchi, M. J. Wainwright, and M. I. Jordan. Local privacy and minimax bounds: Sharp rates for probability estimation. In *NIPS*, pages 1529–1537, 2013.
- [12] C. Dwork, F. McSherry, K. Nissim, and A. Smith. Calibrating noise to sensitivity in private data analysis. In *TCC*, pages 265–284, 2006.
- [13] C. Dwork, A. Roth, et al. The algorithmic foundations of differential privacy. *Foundations and Trends® in Theoretical Computer Science*, 9(3–4):211–407, 2014.
- [14] Ú. Erlingsson, V. Pihur, and A. Korolova. RAPPOR: randomized aggregatable privacy-preserving ordinal response. In *CCS*, pages 1054–1067, 2014.
- [15] A. Evfimievski, J. Gehrke, and R. Srikant. Limiting privacy breaches in privacy preserving data mining. In *PODS*, pages 211–222, 2003.
- [16] G. C. Fanti, V. Pihur, and Ú. Erlingsson. Building a RAPPOR with the unknown: Privacy-preserving learning of associations and data dictionaries. *PoPETs*, 2016(3):41–61, 2016.
- [17] J. Hsu, S. Khanna, and A. Roth. Distributed private heavy hitters. In *Automata, Languages, and Programming - 39th International Colloquium, ICALP 2012, Warwick, UK, July 9-13, 2012, Proceedings, Part I*, pages 461–472, 2012.
- [18] X. Hu, M. Yuan, J. Yao, Y. Deng, L. Chen, Q. Yang, H. Guan, and J. Zeng. Differential privacy in telco big data platform. *PVLDB*, 8(12):1692–1703, 2015.
- [19] N. Mishra and M. Sandler. Privacy via pseudorandom sketches. In *Proceedings of the Twenty-fifth ACM SIGMOD-SIGACT-SIGART Symposium on Principles of Database Systems*, PODS ’06, pages 143–152, New York, NY, USA, 2006. ACM.
- [20] A. Thakurta, A. Vyrros, U. Vaishampayan, A. Kapoor, J. Freudiger, V. Sridhar, and D. Davidson. Learning new words. In *US Patent 9594741*, 2017.

- 196 [21] S. L. Warner. Randomized response: A survey technique for eliminating evasive answer bias.
197 *Journal of the American Statistical Association*, 60(309):63–69, 1965.
- 198 [22] L. Wasserman and S. Zhou. A statistical framework for differential privacy. *Journal of the*
199 *American Statistical Association*, 105(489):375–389, 2010.