# Differentially Private Integer Partitions and their Applications

Jeremiah Blocki

Microsoft Research

### Abstract

Given a positive integer $N \geq 0$ a partition of $N$ is a non-increasing sequence of numbers $x_1 \geq x_2 \ldots \geq x_N \geq 0$ such that $x_1 + \ldots + x_N = N$. We say that two partitions $x$ and $y$ are neighbors if the L1 distance between them is at most $1/2$. Blocki et al. [BDB16] recently showed that there is a $(\epsilon, \delta)$-differentially private algorithm which (whp) achieves L1 error $O(\sqrt{N}/\epsilon)$ and they used their algorithm to publish password frequency data from a password dataset of 70 million Yahoo! users [Bon12]. Their algorithm, which was based on an (approximate) instantiation of the exponential mechanism [MT07], is computationally efficient in $\frac{1}{\epsilon}$, $\log\left(\frac{1}{\delta}\right)$ and $N$.

The applications of the mechanism of Blocki et al. [BDB16] are not limited to passwords. For example, the degree distribution of a social network $G$ is simply a partition of the integer $2\left|E(G)\right|$. Thus, the mechanism of [BDB16] could be used to preserve differential privacy when releasing the degree distribution. It is particularly important to understand the performance of the exponential mechanism for integer partitions. We provide a pure $\epsilon$-differentially instantiation of the exponential mechanism whenever there is an a priori known upper bound on $N$. We also upper bound the mean squared error of the exponential mechanism $O\left(\frac{\sqrt{N}\log^2 N}{\epsilon^2}\right)$. For comparison, the best known results, due to Hay et al. [HLMJ09], achieved mean squared error $O\left(\frac{\sqrt{N}\log^3 N}{\epsilon^2}\right)$.

Additionally, we conjecture that the L1 error of the exponential mechanism scales with $1/\sqrt{\epsilon}$ instead of $1/\epsilon$. Empirical data from the RockYou password frequency dataset supports this conjecture. The conjecture, if true, could lead to the development of several useful node-differentially private algorithms.

## 1 Introduction

A partition of a non-negative integer $n \in \mathbb{N}$ is an ordered list of $n$ integers $x_1 \geq x_2 \geq \ldots \geq x_n \geq 0$ such that $\sum_{i=1}^{n} x_i = n$. We use $\mathcal{P}(n)$ to denote the set of all partitions of the integer $n$. For example, $\mathcal{P}(3) = \{(1,1,1), (2,1,0), (3,0,0)\}$. We let $\mathcal{P} \doteq \bigcup_{n=0}^{\infty} \mathcal{P}(n)$ denote the set of all integer partitions. We say that two partitions $x \in \mathcal{P}$ and $y \in \mathcal{P}$ are adjacent if $\|x - y\|_1 = \sum_{i=1}^{\infty} |x_i - y_i| \leq \frac{1}{2}$ — note that if $x \in \mathcal{P}(n_x)$ then we define $x_i = 0$ for $i > n_x$. A randomized algorithm $A : \mathcal{P} \to \mathcal{P}$ satisfies $(\epsilon, \delta)$-differential privacy if for every pair of adjacent partitions $x, y \in \mathcal{P}$ and every subset $S \subseteq \mathcal{P}$ we have $\Pr[A(x) \in S] \leq e^{\epsilon} \Pr[A(y) \in S] + \delta$.

Given an input dataset $f \in \mathcal{P}$ the exponential mechanism [MT07] $\mathcal{E}^{\epsilon}(f)$ simply outputs each possible outcome $\tilde{f} \in \mathcal{P}$ with probability proportional to $\exp\left(-\frac{\epsilon \cdot \|f - \tilde{f}\|_1}{2}\right)$. While there are infinitely many integer partitions, Blocki et al. [BDB16] observed that this distribution is well defined[1] The exponential mechanism of McSherry and Talwar [MT07] is known to have many powerful applications, especially in regards to the development of differentially private mechanisms in the non-interactive setting [BLR13]. Unfortunately, there is no efficient sampling algorithm instantiating this mechanism in general. Indeed, there is powerful evidence (e.g., [Ull13]) that no general-purpose instantiation of the exponential mechanism can be computationally efficient. However, Blocki et al. [BDB16] recently showed that, for the restricted case of integer partitions, there is an efficient algorithm that (approximately) samples from the exponential mechanism over the space of all integer partitions. More precisely, their algorithm preserves $(\epsilon, \delta)$-differential privacy and runs in

---

[1]Intuitively, if $f \in \mathcal{P}(n)$ then any partition $f' \in \mathcal{P}(n + d)$ has $\|f' - f\| \geq d$. Hardy and Ramanujan [HR18] showed that $|\mathcal{P}(n)| \sim \frac{1}{4n\sqrt{3}} \exp\left(\pi\sqrt{\frac{2n}{3}}\right)$. Thus, the sum $\sum_{d=0} e^{-d\epsilon} |\mathcal{P}(n + d)|$ converges.

polynomial time in $N$, $1/\epsilon$ and $\log \delta$. They also showed that the expected L1 error of their mechanism was upper bounded by $O\left(\frac{\sqrt{N}}{\epsilon}\right)$.

## 1.1 Applications

A password frequency list is simply partition of $f_1 + \ldots + f_N = N$, the number of of user passwords. Password frequency lists from empirical datasets have great value to security researchers who wish to understand the nature of an underlying password distribution so that they can accurately estimate security risks or evaluate various password defenses. For example, the sum $\lambda_\beta = \sum_{i=1}^{\beta} f_i$ is an approximate upper bound on the number of accounts that an untargeted adversary could compromise with $\beta$ guesses per user. Despite their usefulness an organization may understandably be wary of publishing password frequency lists for its own users due to potential security and privacy risks. For example, Yahoo! allowed Bonneau [Bon12] to collect anonymized password frequency data from a random sample of 70 million users and publish some aggregate statistics such as min-entropy. However, Yahoo! declined to publish the original password frequency lists so that other researchers could use them.

Recently, Yahoo! gave Blocki et al. permission to run their differentially private algorithm and place the sanitized password frequency data in the public domain and made it freely available for download.[2]. This is an excellent instance of differential privacy *enabling* (instead of hindering) research by alleviating potential security and privacy concerns. The data has since been used to analyze password hashing algorithms [BD16].

While Blocki et al. focused on password frequency lists, the differentially private algorithm for releasing integer partitions may also be useful in other settings. The degree distribution of a graph $G$ with $n$ nodes and $m$ edges is simply a partition of the number $2m$. Previous differentially private research has focused on releasing the degree distribution of $G$ in both the edge-adjacency [HLMJ09] and vertex adjacency models[KNRS13, RS15].

## 2 New Results

In this section we briefly overview (without proof) several new results and discuss their potential applications.

**Pure-DP.** Assume that we are given (a priori) an upper bound $\hat{N}$ on the integer $N$. For example, if we know that a social network $G \in \mathcal{G}^n$ has $n$ nodes then we know that the social network has at most $m \leq \binom{n}{2}$ edges and therefore the degree distribution is a partition of an integer $N \leq \hat{N} \leq 2\binom{n}{2}$. In this case it makes sense to define the exponential mechanism $\mathcal{E}^\epsilon : \mathcal{P}\left(\leq \hat{N}\right) \to \mathcal{P}\left(\leq \hat{N}\right)$ over the restricted set of partitions $\mathcal{P}\left(\leq \hat{N}\right)$. As it turns out there is a polynomial time algorithm (in $\hat{N}, 1/\epsilon$) to (exactly) sample from this distribution.

While the running time of this new algorithm is worse than the algorithm of Blocki et al. [BDB16], the new algorithm preserves pure $\epsilon$-differential privacy. Thus, the new algorithm can ensure significantly better group privacy guarantees making it a useful tool for achieving node-level differential privacy[3]

**Mean Squared Error.** In their analysis of the exponential mechanism Blocki et al. [BDB16] focused on upper bounding expected L1 error. Their upper bound $O\left(\sqrt{N}/\epsilon\right)$ on L1 error was incomparable to a result of Hay et al. [HLMJ09], who presented an algorithm that achieves mean squared error $O\left(\frac{\sqrt{N}\log^3 N}{\epsilon^2}\right)$. We

---

[2]https://figshare.com/articles/Yahoo_Password_Frequency_Corpus/2057937    main   data   file   SHA-256   hash: 061137ea3cc129c7d9f501295cb194e0c6fa158acac702f893cba3cfd5f44efe

[3]One tempting way to achieve node differential privacy might be to set $\epsilon' = \epsilon/n$ and $\delta' = \delta/n$ and run a $(\epsilon', \delta')$-edge differentially private algorithm. However, this algorithm would only guarantee $(\epsilon, \delta e^{\epsilon n}/n)$-node level differential privacy. This example illustrates why it is desirable to have pure $\epsilon$-differentially private implementation of the exponential mechanism for integer partitions.
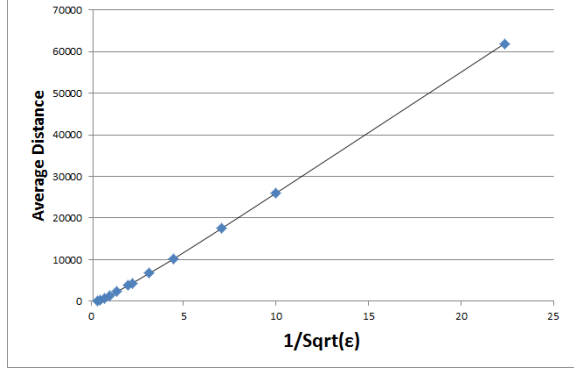
Figure 1: RockYou Password Frequency Data $N \approx 3.26 \times 10^7$. Average L1 error vs $\frac{1}{\sqrt{\epsilon}}$. L1 error is averaged over 100 independent samples.

can now upper bound the mean squared error of the exponential mechanism — $O\left(\frac{\sqrt{N}\log^2 N}{\epsilon^2}\right)$. Thus, the exponential mechanism also improves on L2 error.

**Lower Bound.** We can show that *any* differentially private algorithm for releasing integer partitions must incur L1 error at least $\Omega\left(\frac{\sqrt{N}}{\log N}\right)$ in the worst case. Thus, the exponential mechanism is nearly optimal.

## 3 An Open Question

Blocki et al. [BDB16] proved that if $1/\epsilon = o(\sqrt{N})$ then (whp) the L1 error of the exponential mechanism is at most $c\sqrt{N}/\epsilon$ for some constant $c$. Conjecture 1 says that L1 error of the exponential mechanism actually scales with $\frac{1}{\sqrt{\epsilon}}$ instead of $\frac{1}{\epsilon}$.

**Conjecture 1** *Let* $f \in \mathcal{P}(N)$, $\epsilon > \frac{48\pi^2}{\sqrt{N}}$ *and* $\delta \geq e^{1-\sqrt{N}/2}$ *be given be given and let* $\tilde{f} \leftarrow \mathcal{E}^\epsilon(f)$ *denote a random sample from the exponential mechanism. Then except with probability* $\delta$ *we will have*

$$\|f - \tilde{f}\|_1 \leq \frac{c_1\sqrt{N}}{\sqrt{\epsilon}} + \frac{c_2 \ln\left(\frac{1}{\delta}\right)}{\epsilon} \ ,$$

*where* $c_1$ *and* $c_2$ *are constants.*

### 3.1 Empirical Evidence

Empirical results of Blocki et al. [BDB16] support Conjecture 1. In particular, Blocki et al. [BDB16] analyzed the performance of the exponential mechanism on the RockYou password frequency dataset ($N \approx 3.26 \times 10^7$). Figure 1 plots $\|f - \tilde{f}\|_1$ versus $1/\sqrt{\epsilon}$. Here, $f$ represents the original RockYou frequency data (e.g., $f_1 \approx 3 \times 10^5$ denotes the number of RockYou users who selected the most popular password '123456') and $\tilde{f} \leftarrow \mathcal{E}^\epsilon(f)$ represents a sample from the exponential mechanism. Each point represents the average value $\|f - \tilde{f}\|_1$ taken over 100 independent samples from the exponential mechanism. While the results from Figure 1 certainly do not constitute a proof of Conjecture 1 they are highly suggestive.

### 3.2 Implications for Node Privacy in Social Networks

There are two variants of differential privacy for social networks: edge privacy and node privacy. Intuitively, edge privacy protects each individual link (e.g., relationships), while the later protects an individual together

3

with all of the edges incident to that individual (e.g., all of his/her relationships). Edge privacy, the weaker notion, has been studied more extensively as it is often easier to obtain positive results in this setting (e.g., [BBDS12, HLMJ09, NRS07, KRSY14, KS12]). However, the guarantee of edge privacy may not be a sufficient privacy guarantee in many contexts[4].

While node differential privacy provides much better privacy guarantees it is much harder to develop differentially private algorithms that give accurate answers (e.g., even the simple query how many edges exist in the graph has sensitivity $O(n)$ because we can destroy $n-1$ edges by deleting a single node). Recently, Blocki et al., Kasiviswanathan et al. and Chen and Zhou al. [BBDS13, KNRS13, CZ13] began developing techniques (e.g., lipshitz extensions) for building node differentially private algorithms which would give accurate answers whenever the underling graph satisfied certain sparsity conditions. While this progress has been exciting, these works considered only a very limited classes of queries like subgraph counting queries. Raskhodnikova and Smith [RS15] recently explored the possibility of releasing the degree-distribution with node-differential privacy. However, Raskhodnikova and Smith only promise accurate results when the graph satisfies a stronger assumptions (e.g. $\alpha$-decay).

If Conjecture 1 holds then we can promise accurate results under more general conditions — Theorem 1.
**Notation:** Let $\textbf{DegList}(G)$ denote the degree distribution of a social network $G \in \mathcal{G}^n$, the set of social networks on $n$ nodes. Let $\mathcal{G}_d^n$ denote the set of all social networks on with maximum degree $d$, and let $\mathcal{G}_{d,k}^n$ denote the set of social networks $G$ that are $k$-close to some social network $G'$ with maximum degree $d$ — formally $G$ is $k$-close to $\mathcal{G}_d^n$ if $\min_{G' \in \mathcal{G}_d^n} \|\textbf{DegList}(G) - \textbf{DegList}(G')\|_1 \leq m$.

**Theorem 1** *Assuming that Conjecture 1 holds there is an efficient, in $1/\epsilon, n$, $\epsilon$-differentially private algorithm A such that for any graph $G \in \mathcal{G}_{d,k}^n$*

$$\mathbf{E}\left[\|f - A(f)\|_1\right] = O\left(k + \frac{\sqrt{md}}{\sqrt{\epsilon}}\right) \ .$$

*Here, $f \in \mathcal{P}$ denotes the true degree distribution of $G$. In particular, if $md = o(n^2)$ then $\mathbf{E}\left[\|f - A(f)\|_1\right] = o(n)$.*

*Proof of Theorem 1.* (sketch) [RS15] showed that there is a polynomial time computable 'Lipshitz extension' of the function $\textbf{DegList} : \mathcal{G}^n \to \mathbb{R}^n$. Specifically, there is an efficiently computable function $g_d : \mathcal{G}^n \to \mathbb{R}^n$ with the following properties: (1) $g_d(G) = \textbf{DegList}(G)$ for all $G \in \mathcal{G}_d^n$, (2) for all node-adjacent social networks $G \sim G'$ we have $\|\textbf{DegList}(G) - \textbf{DegList}(G')\|_1 \leq 3d$, and (3) $\|\textbf{DegList}(G) - g_d(G) - \|_1 \leq O(k)$ for all $G \in \mathcal{G}_{d,k}^n$. Thus, we simple output $\mathcal{E}^{\epsilon/(3d)}\left(g_d(G)\right)$. The mechanism preserves $\epsilon$-differential privacy because the global sensitivity of $g_d$ is at most $3d$ (property 2). By Conjecture 1 we have

$$\|g_d(G) - \mathcal{E}^{\epsilon/(3d)}\left(g_d(G)\right)\|_1 \leq O\left(\frac{\sqrt{dm}}{\sqrt{\epsilon}} - \frac{\ln \delta}{\epsilon}\right) \ ,$$

except with probability $\delta$. For $G \in \mathcal{G}_{d,k}^n$ we can apply the triangle inequality and (3) to obtain

$$\|\textbf{DegList}(G) - \mathcal{E}^{\epsilon/(3d)}\left(g_d(G)\right)\|_1 \leq O(k) + O\left(\frac{\sqrt{dm}}{\sqrt{\epsilon}} - \frac{\ln \delta}{\epsilon}\right) \ .$$

$\square$

---

[4]For example consider the problem of releasing statistics about the communication graph for a site like Ashley Madison, where an adversary might have significant background knowledge about users due to the infamous data breach. While edge privacy would not disclose the existence of particular relationship, it would not hide a user's overall activity level.

# References

[BBDS12] Jeremiah Blocki, Avrim Blum, Anupam Datta, and Or Sheffet. The johnson-lindenstrauss transform itself preserves differential privacy. In *Foundations of Computer Science (FOCS), 2012 IEEE 53rd Annual Symposium on*, pages 410–419. IEEE, 2012.

[BBDS13] Jeremiah Blocki, Avrim Blum, Anupam Datta, and Or Sheffet. Differentially private data analysis of social networks via restricted sensitivity. In *Proceedings of the 4th conference on Innovations in Theoretical Computer Science*, pages 87–96. ACM, 2013.

[BD16] Jeremiah Blocki and Anupam Datta. Cash: A cost asymmetric secure hash algorithm for optimal password protection. In *Computer Security Foundations Symposium (CSF), 2016 IEEE 29th*, page (to appear). IEEE, 2016.

[BDB16] Jeremiah Blocki, Anupam Datta, and Joseph Bonneau. Differentially private password frequency lists: Or, how to release statistics from 70 million passwords (on purpose). In *23rd Annual Network and Distributed System Security Symposium, NDSS 2016, San Diego, California, USA, February 21-24, 2016*, 2016.

[BLR13] Avrim Blum, Katrina Ligett, and Aaron Roth. A learning theory approach to noninteractive database privacy. *Journal of the ACM (JACM)*, 60(2):12, 2013.

[Bon12] J. Bonneau. The science of guessing: analyzing an anonymized corpus of 70 million passwords. In *Security and Privacy (SP), 2012 IEEE Symposium on*, pages 538–552. IEEE, 2012.

[CZ13] Shixi Chen and Shuigeng Zhou. Recursive mechanism: towards node differential privacy and unrestricted joins. In *Proceedings of the 2013 ACM SIGMOD International Conference on Management of Data*, pages 653–664. ACM, 2013.

[HLMJ09] Michael Hay, Chao Li, Gerome Miklau, and David Jensen. Accurate estimation of the degree distribution of private networks. In *Data Mining, 2009. ICDM'09. Ninth IEEE International Conference on*, pages 169–178. IEEE, 2009.

[HR18] Godfrey H Hardy and Srinivasa Ramanujan. Asymptotic formulæ in combinatory analysis. *Proceedings of the London Mathematical Society*, 2(1):75–115, 1918.

[KNRS13] Shiva Prasad Kasiviswanathan, Kobbi Nissim, Sofya Raskhodnikova, and Adam Smith. Analyzing graphs with node differential privacy. In *Theory of Cryptography*, pages 457–476. Springer, 2013.

[KRSY14] Vishesh Karwa, Sofya Raskhodnikova, Adam Smith, and Grigory Yaroslavtsev. Private analysis of graph structure. *ACM Transactions on Database Systems (TODS)*, 39(3):22, 2014.

[KS12] Vishesh Karwa and Aleksandra B Slavković. Differentially private graphical degree sequences and synthetic graphs. In *Privacy in Statistical Databases*, pages 273–285. Springer, 2012.

[MT07] Frank McSherry and Kunal Talwar. Mechanism design via differential privacy. In *Foundations of Computer Science, 2007. FOCS'07. 48th Annual IEEE Symposium on*, pages 94–103. IEEE, 2007.

[NRS07] Kobbi Nissim, Sofya Raskhodnikova, and Adam Smith. Smooth sensitivity and sampling in private data analysis. In *Proceedings of the thirty-ninth annual ACM symposium on Theory of computing*, pages 75–84. ACM, 2007.

[RS15] Sofya Raskhodnikova and Adam D. Smith. Efficient lipschitz extensions for high-dimensional graph statistics and node private degree distributions. *CoRR*, abs/1504.07912, 2015.

[Ull13]      Jonathan Ullman. Answering n {2+ o (1)} counting queries with differential privacy is hard. In *Proceedings of the forty-fifth annual ACM symposium on Theory of computing*, pages 361–370. ACM, 2013.