
Discrete Distribution Estimation under Local Privacy

Peter Kairouz *†

Keith Bonawitz *

Daniel Ramage *

KAIROUZ2@ILLINOIS.EDU

BONAWITZ@GOOGLE.COM

DRAMAGE@GOOGLE.COM

* Google, 1600 Amphitheatre Parkway, Mountain View, CA 94043,

† University of Illinois, Urbana-Champaign, 1308 W Main St, Urbana, IL 61801

Abstract

The collection and analysis of user data drives improvements in the app and web ecosystems, but comes with risks to privacy. This paper examines discrete distribution estimation under local privacy, a setting wherein service providers can learn the distribution of a categorical statistic of interest without collecting the underlying data. We present new mechanisms, including hashed k -ary Randomized Response (k -RR), that empirically meet or exceed the utility of existing mechanisms at all privacy levels. New theoretical results demonstrate the order-optimality of k -RR and the existing RAPPOR mechanism at different privacy regimes.

1. Introduction

Software and service providers increasingly see the collection and analysis of user data as key to improving their services. Datasets of user interactions give insight to analysts and provide training data for machine learning models. But the collection of these datasets comes with risk—can the service provider keep the data secure from unauthorized access? Misuse of data can violate the privacy of users and substantially tarnish the provider’s reputation.

One way to minimize risk is to store less data: providers can methodically consider what data to collect and how long to store it. However, even a carefully processed dataset can compromise user privacy. In a now famous study, (Narayanan & Shmatikov, 2008) showed how to de-anonymize watch histories released in the Netflix Prize, a public recommender system competition. While most providers do not intentionally release anonymized datasets, security breaches can mean that even internal, anonymized

datasets have the potential to become privacy problems.

Fortunately, mathematical formulations exist that can give the benefits of population-level statistics without the collection of raw data. Local differential privacy (Duchi et al., 2013a;b) is one such formulation, requiring each device (or session for a cloud service) to share only a noised version of its raw data with the service provider’s logging mechanism. No matter what computation is done to the noised output of a locally differentially private mechanism, any attempt to impute properties of a single record will have a significant probability of error. But not all differentially private mechanisms are equal when it comes to utility: some mechanisms have better accuracy than others for a given analysis, amount of data, and desired privacy level.

Private distribution estimation. This paper investigates the fundamental problem of discrete distribution estimation under local differential privacy. We focus on discrete distribution estimation because it enables a variety of useful capabilities, including usage statistics breakdowns and count-based machine learning models, e.g. naive Bayes (McCallum et al., 1998). We consider empirical, maximum likelihood, and minimax distribution estimation, and study the price of local differential privacy under a variety of loss functions and privacy regimes. In particular, we compare the performance of two recent local privacy mechanisms: (a) the Randomized Aggregatable Privacy-Preserving Ordinal Response (RAPPOR) (Erlingsson et al., 2014), and (b) the k -ary Randomized Response (k -RR) (Kairouz et al., 2014) from a theoretical and empirical perspective.

Our contributions are:

1. For binary alphabets, we prove that Warner’s randomized response model (Warner, 1965) is globally optimal for any loss function and any privacy level (Section 3).
2. For k -ary alphabets, we show that RAPPOR is order optimal in the high privacy regime and strictly sub-optimal in the low privacy regime for ℓ_1 and ℓ_2 losses using an empirical estimator. Conversely, k -RR is order optimal in the low privacy regime and strictly sub-optimal in the

high privacy regime (Section 4.1).

3. Large scale simulations show that the optimal decoding algorithm for both k -RR and RAPPOR depends on the shape of the true underlying distribution. For skewed distributions, the *projected estimator* (introduced here) offers the best utility across a wide variety of privacy levels and sample sizes (Section 4.4).
4. For open alphabets in which the set of input symbols is not enumerable *a priori* we construct the O-RR mechanism (an extension to k -RR using hash functions and cohorts) and provide empirical evidence that the performance of O-RR meets or exceeds that of RAPPOR over a wide range of privacy settings (Section 5).
5. We apply the O-RR mechanism to closed k -ary alphabets, replacing hash functions with permutations. We provide empirical evidence that the performance of O-RR meets or exceeds that of k -RR and RAPPOR in both low and high privacy regimes (Section 5.4).

Related work. There is a rich literature on distribution estimation under local privacy (Chan et al., 2012; Hsu et al., 2012; Bassily & Smith, 2015), of which several works are particularly relevant herein. (Warner, 1965) was the first to study the local privacy setting and propose the randomized response model that will be detailed in Section 3. (Kairouz et al., 2014) introduced k -RR and showed that it is optimal in the low privacy regime for a rich class of information theoretic utility functions. k -RR will be extended to open alphabets in Section 5.1. (Duchi et al., 2013a;b) was the first to apply differential privacy to the local setting, to study the fundamental trade-off between privacy and minimax distribution estimation in the high privacy regime, and to introduce the core of k -RAPPOR. (Erlingsson et al., 2014) proposed RAPPOR, systematically addressing a variety of practical issues for private distribution estimation, including robustness to attackers with access to multiple reports over time, and estimating distributions over open alphabets. RAPPOR has been deployed in the Chrome browser to allow Google to privately monitor the impact of malware on homepage settings. RAPPOR will be investigated in Sections 4.2 and 5.2.

Private distribution estimation also appears in the global privacy context where a trusted service provider releases randomized data (e.g., NIH releasing medical records) to protect sensitive user information (Dwork, 2006; Dwork et al., 2006; Dwork & Lei, 2009; Dwork, 2008; Diakonikolas et al., 2015; Blocki et al., 2016).

2. Preliminaries

2.1. Local differential privacy

Let X be a private source of information defined on a discrete, finite input alphabet $\mathcal{X} = \{x_1, \dots, x_k\}$. A statistical privatization mechanism is a family of distributions \mathcal{Q} that map $X = x$ to $Y = y$ with probability $\mathcal{Q}(y|x)$. Y , the privatized version of X , is defined on an output alphabet $\mathcal{Y} = \{y_1, \dots, y_l\}$ that need not be identical to the input alphabet \mathcal{X} . In this paper, we will represent a privatization mechanism \mathcal{Q} via a $k \times l$ row-stochastic matrix. A conditional distribution \mathcal{Q} is said to be ε -locally differentially private if for all $x, x' \in \mathcal{X}$ and all $E \subset \mathcal{Y}$, we have that

$$\mathcal{Q}(E|x) \leq e^\varepsilon \mathcal{Q}(E|x'), \quad (1)$$

where $\mathcal{Q}(E|x) = \mathbb{P}(Y \in E|X = x)$ and $\varepsilon \in [0, \infty)$ (Duchi et al., 2013a). In other words, by observing $Y \in E$, the adversary cannot reliably infer whether $X = x$ or $X = x'$ (for any pair x and x'). Indeed, the smaller the ε is, the closer the likelihood ratio of $X = x$ to $X = x'$ is to 1. Therefore, when ε is small, the adversary cannot recover the true value of X reliably.

2.2. Private distribution estimation

The private multinomial estimation problem is defined as follows. Given a vector $\mathbf{p} = (p_1, \dots, p_k)$ on the probability simplex \mathbb{S}^k , samples X_1, \dots, X_n are drawn i.i.d. according to \mathbf{p} . An ε -locally differentially private mechanism \mathcal{Q} is then applied independently to each sample X_i to produce $Y^n = (Y_1, \dots, Y_n)$, the sequence of private observations. Observe that the Y_i 's are distributed according to $\mathbf{m} = \mathbf{p}\mathcal{Q}$ and not \mathbf{p} . Our goal is to estimate the distribution vector \mathbf{p} from Y^n .

Privacy vs. utility. There is a fundamental trade-off between utility and privacy. The more private you want to be, the less utility you can get. To formally analyze the privacy-utility trade-off, we study the following constrained minimization problem

$$r_{\ell, \varepsilon, k, n} = \inf_{\mathcal{Q} \in \mathcal{D}_\varepsilon} r_{\ell, \varepsilon, k, n}(\mathcal{Q}), \quad (2)$$

where

$$r_{\ell, \varepsilon, k, n}(\mathcal{Q}) = \inf_{\hat{\mathbf{p}}} \sup_{\mathbf{p}} \mathbb{E}_{Y^n \sim \mathbf{p}\mathcal{Q}} \ell(\mathbf{p}, \hat{\mathbf{p}})$$

is the minimax risk under \mathcal{Q} , ℓ is an application dependent loss function, and \mathcal{D}_ε is the set of all ε -locally differentially private mechanisms.

This problem, though of great value, is intractable in general. Indeed, finding minimax estimators in the non-private setting is already hard for several loss functions. For instance, the minimax estimator under ℓ_1 loss is unknown

even until today. However, in the high privacy regime, we are able to bound the minimax risk of any differentially private mechanism \mathcal{Q} .

Proposition 1 *For the private distribution estimation problem in (2), for any ε -locally differentially private mechanism \mathcal{Q} , there exist universal constants $0 < c_l \leq c_u < 5$ such that for all $\varepsilon \in [0, 1]$,*

$$c_l \min \left\{ 1, \frac{1}{\sqrt{n\varepsilon^2}}, \frac{k}{n\varepsilon^2} \right\} \leq r_{\ell_2, \varepsilon, k, n} \leq c_u \min \left\{ 1, \frac{k}{n\varepsilon^2} \right\},$$

and

$$c_l \min \left\{ 1, \frac{k}{\sqrt{n\varepsilon^2}} \right\} \leq r_{\ell_1, \varepsilon, k, n} \leq c_u \min \left\{ 1, \frac{k}{\sqrt{n\varepsilon^2}} \right\}$$

Proof See (Duchi et al., 2013b). \blacksquare

This result shows that in the high privacy regime ($\varepsilon \leq 1$), the effective sample size of a dataset decreases from n to $n\varepsilon^2/k$. In other words, a factor of k/ε^2 extra samples are needed to achieve the same minimax risk. This is problematic for large alphabets. Our work shows that (a) this problem can be (partially) circumvented using a combination of cohort-style hashing and k -RR (Section 5), and (b) the dependence on the alphabet size vanishes in the moderate to low privacy regime (Section 4.3).

3. Binary Alphabets

In this section, we study the problem of private distribution estimation under binary alphabets. In particular, we show that Warner’s randomized response model (W-RR) is optimal for binary distribution minimax estimation (Warner, 1965). In W-RR, interviewees flip a biased coin (that only they can see the result of), such that a fraction η of participants answer the question “Is the predicate P true (of you)?” while the remaining participants answer the negation (“Is $\neg P$ true?”), without revealing which question they answered. For $\eta = e^\varepsilon$ ($\varepsilon \geq 0$), W-RR can be described by the following 2×2 row-stochastic matrix

$$\mathcal{Q}_{\text{WRR}} = \frac{1}{e^\varepsilon + 1} \begin{bmatrix} e^\varepsilon & 1 \\ 1 & e^\varepsilon \end{bmatrix}. \quad (3)$$

It is easy to check that the above mechanism satisfies the constraints imposed by local differential privacy.

Theorem 2 *For all binary distributions \mathbf{p} , all loss functions ℓ , and all privacy levels ε , \mathcal{Q}_{WRR} is the optimal solution to the private minimax distribution estimation problem in (2).*

Proof sketch. (Kairouz et al., 2014) showed that W-RR dominates all other differentially private mechanisms in a

strong Markovian sense: for any binary differentially private mechanism \mathcal{Q} , there exists a 2×2 stochastic mapping \mathcal{W} such that $\mathcal{Q} = \mathcal{W} \circ \mathcal{Q}_{\text{WRR}}$. Therefore, for any risk function $r(\cdot)$ that obeys the data processing inequality ($r(\mathcal{Q}) \leq r(\mathcal{Q} \circ \mathcal{W})$ for any stochastic mappings \mathcal{Q} and \mathcal{W}), we have that $r(\mathcal{Q}_{\text{WRR}}) \leq r(\mathcal{Q})$ for any binary differentially private mechanism \mathcal{Q} . In Supplementary Section A, we prove that $r_{\ell, \varepsilon, k, n}(\mathcal{Q})$ obeys the data processing inequality, thus W-RR achieves the optimal privacy-utility trade-off under minimax distribution estimation.

4. k -ary Alphabets

Above, we saw that W-RR is optimal for all privacy levels and all loss functions. However, it can only be applied to binary alphabets. In this section, we study optimal privacy mechanisms for k -ary alphabets. We show that under ℓ_1 and ℓ_2 losses, k -RAPPOR is order optimal in the high privacy regime and sub-optimal in the low privacy regime. Conversely, k -RR is order optimal in the low privacy regime and sub-optimal in the high privacy regime.

4.1. The k -ary Randomized Response

The k -ary randomized response (k -RR) mechanism is a locally differentially private mechanism that maps \mathcal{X} stochastically onto itself (i.e., $\mathcal{Y} = \mathcal{X}$), given by

$$\mathcal{Q}_{\text{KRR}}(y|x) = \frac{1}{k-1+e^\varepsilon} \begin{cases} e^\varepsilon & \text{if } y = x, \\ 1 & \text{if } y \neq x. \end{cases} \quad (4)$$

k -RR can be viewed as a multiple choice generalization of the W-RR mechanism (note that k -RR reduces to W-RR for $k = 2$). In (Kairouz et al., 2014), the k -RR mechanism was shown to be optimal in the low privacy regime for a large class of information theoretic utility functions.

Empirical estimation under k -RR. It is easy to see that under \mathcal{Q}_{KRR} , outputs are distributed according to:

$$\mathbf{m} = \frac{e^\varepsilon - 1}{e^\varepsilon + k - 1} \mathbf{p} + \frac{1}{e^\varepsilon + k - 1} \quad (5)$$

The empirical estimate of \mathbf{p} under \mathcal{Q}_{KRR} is given by

$$\begin{aligned} \hat{\mathbf{p}} &= \hat{\mathbf{m}} \mathcal{Q}_{\text{KRR}}^{-1} \\ &= \frac{e^\varepsilon + k - 1}{e^\varepsilon - 1} \hat{\mathbf{m}} - \frac{1}{e^\varepsilon - 1}, \end{aligned} \quad (6)$$

where $\hat{\mathbf{m}}$ is the empirical estimate of \mathbf{m} and

$$\mathcal{Q}_{\text{KRR}}^{-1}(y|x) = \frac{1}{e^\varepsilon - 1} \begin{cases} e^\varepsilon + k - 2 & \text{if } y = x, \\ -1 & \text{if } y \neq x. \end{cases} \quad (7)$$

via the Sherman-Morrison formula. Observe that because $\hat{\mathbf{m}} \rightarrow \mathbf{m}$ almost surely, $\hat{\mathbf{p}} \rightarrow \mathbf{p}$ almost surely.

Proposition 3 For the private distribution estimation problem under k -RR and its empirical estimator given in (6), for all ε , n , and k , we have that

$$\mathbb{E} \ell_2^2(\hat{\mathbf{p}}, \mathbf{p}) = \frac{1 - \sum_{i=1}^k p_i^2}{n} + \frac{k-1}{n} \left(\frac{k + 2(e^\varepsilon - 1)}{(e^\varepsilon - 1)^2} \right),$$

and for large n , $\mathbb{E} \ell_1(\hat{\mathbf{p}}, \mathbf{p}) \approx$

$$\sum_{i=1}^k \sqrt{\frac{2((e^\varepsilon - 1)p_i + 1)((e^\varepsilon - 1)(1 - p_i) + k - 1)}{\pi n (e^\varepsilon - 1)^2}},$$

where $a_n \approx b_n$ means $\lim_{n \rightarrow \infty} a_n/b_n = 1$.

Proof See Supplementary Section B. ■

Observe that for $\mathbf{p}_U = (\frac{1}{k}, \dots, \frac{1}{k})$, we have that

$$\begin{aligned} \mathbb{E} \ell_2^2(\hat{\mathbf{p}}, \mathbf{p}) &\leq \mathbb{E} \ell_2^2(\hat{\mathbf{p}}, \mathbf{p}_U) \\ &= \left(1 + \frac{k + 2(e^\varepsilon - 1)}{(e^\varepsilon - 1)^2} k \right) \frac{1 - \frac{1}{k}}{n}, \end{aligned} \quad (8)$$

and

$$\begin{aligned} \mathbb{E} \ell_1(\hat{\mathbf{p}}, \mathbf{p}) &\leq \mathbb{E} \ell_1(\hat{\mathbf{p}}, \mathbf{p}_U) \\ &\approx \left(\frac{e^\varepsilon + k - 1}{e^\varepsilon - 1} \right) \sqrt{\frac{2(k-1)}{\pi n}}. \end{aligned} \quad (9)$$

Constraining empirical estimates to \mathbb{S}^k . It is easy to see that $\|\hat{\mathbf{p}}_{\text{KRR}}\|_1 = 1$. However, some of the entries of $\hat{\mathbf{p}}_{\text{KRR}}$ can be negative (especially for small values of n). Several remedies are available, including (a) truncating the negative entries to zero and renormalizing the entire vector to sum to 1, or (b) projecting $\hat{\mathbf{p}}_{\text{KRR}}$ onto the probability simplex. We evaluate both approaches in Section 4.4.

4.2. k -RAPPOR

The randomized aggregatable privacy-preserving ordinal response (RAPPOR) is an open source Google technology for collecting aggregate statistics from end-users with strong local differential privacy guarantees (Erlingsson et al., 2014). The simplest version of RAPPOR, called the basic one-time RAPPOR and referred to herein as k -RAPPOR, first appeared in (Duchi et al., 2013a;b). k -RAPPOR maps the input alphabet \mathcal{X} of size k to an output alphabet \mathcal{Y} of size 2^k . In k -RAPPOR, we first map \mathcal{X} deterministically to $\tilde{\mathcal{X}} = \mathbb{R}^k$, the k -dimensional Euclidean space. Precisely, $X = x_i$ is mapped to $\tilde{X} = e_i$, the i^{th} standard basis vector in \mathbb{R}^k . We then randomize the coordinates of \tilde{X} independently to obtain the private vector $Y \in \{0, 1\}^k$. Formally, the j^{th} coordinate of Y is given by: $Y^{(j)} = \tilde{X}^{(j)}$ with probability $e^{\varepsilon/2}/(1 + e^{\varepsilon/2})$ and $1 - \tilde{X}^{(j)}$ with probability $1/(1 + e^{\varepsilon/2})$. The randomization in $\mathbf{Q}_{k\text{-RAPPOR}}$ is ε -locally differentially private (Duchi et al., 2013a; Erlingsson et al., 2014).

Under k -RAPPOR, $Y_i = [Y_i^{(1)}, \dots, Y_i^{(k)}]$ is a k -dimensional binary vector, which implies that

$$\mathbb{P}(Y_i^{(j)} = 1) = \left(\frac{e^{\varepsilon/2} - 1}{e^{\varepsilon/2} + 1} \right) p_j + \frac{1}{e^{\varepsilon/2} + 1}, \quad (10)$$

for all $i \in \{1, \dots, n\}$ and $j \in \{1, \dots, k\}$.

Empirical estimation under k -RAPPOR. Let Y^n be the $n \times k$ matrix formed by stacking the row vectors Y_1, \dots, Y_n on top of each other. The empirical estimator of \mathbf{p} under k -RAPPOR is:

$$\hat{p}_j = \left(\frac{e^{\varepsilon/2} + 1}{e^{\varepsilon/2} - 1} \right) \frac{T_j}{n} - \frac{1}{e^{\varepsilon/2} - 1}, \quad (11)$$

where $T_j = \sum_{i=1}^n Y_i^{(j)}$. Because T_j/n converges to m_j almost surely, \hat{p}_j converges to p_j almost surely. As with k -RR, we can constrain $\hat{\mathbf{p}}$ to \mathbb{S}^k through truncation and normalization or through projection (described in Section 4.1), both of which will be evaluated in Section 4.4.

Proposition 4 For the private distribution estimation problem under k -RAPPOR and its empirical estimator given in (11), for all ε , n , and k , we have that

$$\mathbb{E} \ell_2^2(\hat{\mathbf{p}}, \mathbf{p}) = \frac{1 - \sum_{i=1}^k p_i^2}{n} + \frac{k e^{\varepsilon/2}}{n (e^{\varepsilon/2} - 1)^2},$$

and for large n , $\mathbb{E} \ell_1(\hat{\mathbf{p}}, \mathbf{p}) \approx$

$$\sum_{i=1}^k \sqrt{\frac{2((e^{\varepsilon/2} - 1)p_i + 1)((e^{\varepsilon/2} - 1)(1 - p_i) + 1)}{\pi n (e^{\varepsilon/2} - 1)^2}},$$

where $a_n \approx b_n$ means $\lim_{n \rightarrow \infty} a_n/b_n = 1$.

Proof See Supplementary Section C. ■

Observe that for $\mathbf{p}_U = (\frac{1}{k}, \dots, \frac{1}{k})$, we have that

$$\begin{aligned} \mathbb{E} \ell_2^2(\hat{\mathbf{p}}, \mathbf{p}) &\leq \mathbb{E} \ell_2^2(\hat{\mathbf{p}}, \mathbf{p}_U) \\ &= \left(1 + \frac{k^2 e^{\varepsilon/2}}{(k-1)(e^{\varepsilon/2} - 1)^2} \right) \frac{1 - \frac{1}{k}}{n}, \end{aligned} \quad (12)$$

and

$$\begin{aligned} \mathbb{E} \ell_1(\hat{\mathbf{p}}, \mathbf{p}) &\leq \mathbb{E} \ell_1(\hat{\mathbf{p}}, \mathbf{p}_U) \\ &\approx \sqrt{\frac{(e^{\varepsilon/2} + k - 1)(e^{\varepsilon/2}(k-1) + 1)}{(e^{\varepsilon/2} - 1)^2(k-1)}} \sqrt{\frac{2(k-1)}{\pi n}}. \end{aligned} \quad (13)$$

4.3. Theoretical Analysis

We now analyze the performance of k -RR and k -RAPPOR relative to maximum likelihood estimation (which is equivalent to empirical estimation) on the non-privatized data

X^n . In the non-private setting, the maximum likelihood estimator has a worst case risk of $\sqrt{\frac{2(k-1)}{\pi n}}$ under the ℓ_1 loss, and a worst case risk of $\frac{1-\frac{1}{k}}{n}$ under the ℓ_2^2 loss (Lehmann & Casella, 1998; Kamath et al., 2015).

Performance under k -RR. Comparing Equation (8) to the observation above, we can see that an extra factor of $\left(1 + \frac{k+2(e^\varepsilon-1)}{(e^\varepsilon-1)^2}k\right)$ samples is needed to achieve the same ℓ_2^2 loss as in the non-private setting. Similarly, from Equation (9), a factor of $\left(\frac{e^\varepsilon+k-1}{e^\varepsilon-1}\right)^2$ samples is needed under the ℓ_1 loss. For small ε , the sample size n is effectively reduced to $n\varepsilon^2/k^2$ (under both losses). When compared to Proposition 1, this result implies that k -RR is not optimal in the high privacy regime. However, for $\varepsilon \approx \ln k$, the sample size n is reduced to $n/4$ (under both losses). This result suggests that, while k -RR is not optimal for small values of ε , it is “order” optimal for ε on the order of $\ln k$. Note that k -RR provides a natural interpretation of this low privacy regime: specifically, setting $\varepsilon = \ln k$ translates to telling the truth with probability $\frac{1}{2}$ and lying uniformly over the remainder of the alphabet with probability $\frac{1}{2}$; an intuitively reasonable notion of plausible deniability.

Performance under k -RAPPOR. Comparing Equation (12) to the observation at the beginning of this subsection, we can see that an extra factor of $\left(1 + \frac{k^2 e^{\varepsilon/2}}{(k-1)(e^{\varepsilon/2}-1)^2}\right)$ samples is needed to achieve the same ℓ_2^2 as in the non-private case. Similarly, from Equation (13), an extra factor of $\frac{(e^{\varepsilon/2}+k-1)(e^{\varepsilon/2}(k-1)+1)}{(e^{\varepsilon/2}-1)^2(k-1)}$ samples is needed under the ℓ_1 loss. For small ε , n is effectively reduced to $n\varepsilon^2/4k$ (under both losses). When compared to Proposition 1, this result implies that k -RAPPOR is “order” optimal in the high privacy regime. However, for $\varepsilon \approx \ln k$, n is reduced to n/\sqrt{k} (under both losses). This suggests that k -RAPPOR is strictly sub-optimal in the moderate to low privacy regime.

Proposition 5 For all $\mathbf{p} \in \mathbb{S}^k$ and all $\varepsilon \geq \ln(k/2)$,

$$\mathbb{E} \|\hat{\mathbf{p}}_{KRR} - \mathbf{p}\|_2^2 \leq \mathbb{E} \|\hat{\mathbf{p}}_{RAPPOR} - \mathbf{p}\|_2^2, \quad (14)$$

where $\hat{\mathbf{p}}_{KRR}$ is the empirical estimate of \mathbf{p} under k -RR, $\hat{\mathbf{p}}_{RAPPOR}$ is the empirical estimate of \mathbf{p} under k -RAPPOR, and $\hat{\mathbf{p}}$ is the empirical estimator under k -RAPPOR.

Proof See Supplementary Section D. ■

4.4. Simulation Analysis

To complement the theoretical analysis above, we ran simulations of k -RR and k -RAPPOR varying the alphabet size k , the privacy level ε , the number of users n , and the true distribution \mathbf{p} from which the samples were drawn. In

all cases, we report the mean over 10,000 evaluations of $\|\hat{\mathbf{p}} - \hat{\mathbf{p}}^{\text{decoded}}\|_1$ where $\hat{\mathbf{p}}$ is the ground truth sample drawn from the true distribution and $\hat{\mathbf{p}}^{\text{decoded}}$ is the decoded k -RR or k -RAPPOR distribution. We vary ε over a range that corresponds to the moderate-to-low privacy regimes in our theoretical analysis above, observing that even large values of ε can provide plausible deniability impossible under un-noised logging.

We compare using the ℓ_1 distance of the two distributions because in most applications we want to estimate all values well, emphasizing neither very large values (as an ℓ_2 or higher metric might) nor very small values (as information theoretic metrics might). Supplementary Figures 5 and 6, analogous to the ones in this section, demonstrate that the choice of distance metric does not qualitatively affect our conclusions on the decoding strategies for k -RR or k -RAPPOR nor on the regimes in which each is superior.

The distributions we considered in simulation were binomial distributions with parameter in $\{.1, .2, .3, .4, .5\}$, Zipf distribution with parameter in $\{1, 2, 3, 4, 5\}$, multinomial distributions drawn from a symmetric Dirichlet distribution with parameter $\vec{1}$, and the geometric distribution with mean $k/5$. The geometric distribution is shown in Supplementary Figure 4. We focus primarily on the geometric distribution here because qualitatively it shows the same patterns for decoding as the full set of binomial and Zipf distributions and it is sufficiently skewed to represent many real-world datasets. It is also the distribution for which k -RAPPOR does the best relative to k -RR over the largest range of k and ε in our simulations.

4.4.1. DECODING

We first consider the impact of the choice of decoding mechanism used for k -RR and k -RAPPOR. We find that the best decoder in practice for both k -RR and k -RAPPOR on skewed distributions is the *projected decoder* which projects the $\hat{\mathbf{p}}_{KRR}$ or $\hat{\mathbf{p}}_{RAPPOR}$ onto the probability simplex \mathbb{S}^k using the method described in Algorithm 1 of (Wang & Carreira-Perpiñán, 2013). For k -RR, we compare the projected empirical decoder to the normalized empirical decoder (which truncates negative values and renormalizes) and to the maximum likelihood decoder (see Supplementary Section F.1). For k -RAPPOR, we compare the standard decoder, normalized decoder, and projected decoder. Figure 1 shows that the projected decoder is substantially better than the other decoders for both k -RR and k -RAPPOR for the whole range of k and ε for the geometric distribution. We find this result holds as we vary the number of users from 30 to 10^6 and for all distributions we evaluated except for the Dirichlet distribution, which is the least skewed. For the Dirichlet distribution, the normalized decoder variant is best for both k -RR and k -RAPPOR. Be-

cause the projected decoder is best on all the skewed distributions we expect to see in practice, we use it exclusively for the open-alphabet experiments in Section 5.

4.4.2. k -RR vs k -RAPPOR

To construct a fair, empirical comparison of k -RR and k -RAPPOR, we employ the same methodology used above in selecting decoders. Figure 2 shows the difference between the best k -RR decoder and the best k -RAPPOR decoder (for a particular k and ε). For most cells, the best decoder is the projected decoder described above.

Note that the best k -RAPPOR decoder is consistently better than the best k -RR decoder for relatively large k and low ε . However, k -RR is slightly better than k -RAPPOR in all conditions where $k < e^\varepsilon$ (bottom-right triangle), an empirical result for ℓ_1 that complements Proposition 5’s statement about ML decoders in ℓ_2 . All of the skewed distributions manifest the same pattern as the geometric distribution. As the number of users increases, k -RR’s advantage over k -RAPPOR in the low privacy environment shrinks. In the next sections, we will examine the use of cohorts to improve decoding and to handle larger, open alphabets.

5. Open Alphabets, Hashing, and Cohorts

In practice, the set of values that may need to be collected may not be easily enumerable in advance, preventing a direct application of the binary and k -ary formulations of private distribution estimation. Consider a population of n users, where each user i possesses a symbol s_i drawn from a large set of symbols \mathcal{S} whose membership is not known in advance. This scenario is common in practice; for example, in Chrome’s estimation of the distribution of home page settings (Erlingsson et al., 2014). Building on this intuitive example, we assume for the remainder of the paper that symbols s_i are strings, but we note that the methods described are applicable to any hashable structures.

5.1. O-RR: k -RR with hashing and cohorts

k -RR is effective for privatizing over known alphabets. Inspired by (Erlingsson et al., 2014), we extend k -RR to open alphabets by combining two primary intuitions: hashing and cohorts. Let $\text{HASH}(s)$ be a function mapping $\mathcal{S} \rightarrow \mathbb{N}$ with a low collision rate, i.e. $\text{HASH}(s) = \text{HASH}(s')$ with very low probability for $s' \neq s$. With hashing, we could use k -RR to guarantee local privacy over an alphabet of size k by having each client report $\mathbf{Q}_{\text{KRR}}(\text{HASH}(s) \bmod k)$. However, as we will see, hashing alone is not enough to provide high utility because of the increased rate of collisions introduced by the modulus.

Complementing hashing, we also apply the idea of hash cohorts: each user i is assigned to a cohort c_i sampled i.i.d.

from the uniform distribution over $\mathcal{C} = \{1, \dots, C\}$. Each cohort $c \in \mathcal{C}$ provides an independent view of the underlying distribution of strings by projecting the space of strings \mathcal{S} onto a smaller space of symbols \mathcal{X} using an independent hash function HASH_c . The users in a cohort use their cohort’s hash function to partition \mathcal{S} into k disjoint subsets by computing $x_i = \text{HASH}_{c_i}(s_i) \bmod k = \text{HASH}_{c_i}^{(k)}(s_i)$. Each subset contains approximately the same number of strings, and because each cohort uses a different hash function, the induced partitions for different cohorts are orthogonal: $\mathbb{P}(x_i = x_j | c_i \neq c_j) \approx \frac{1}{k}$ even when $s_i = s_j$.

5.1.1. ENCODING AND DECODING

For encoding, the O-RR privatization mechanism can be viewed as a sampling distribution independent of \mathcal{C} . Therefore, $\mathbf{Q}_{\text{ORR}}(y, c|s)$ is given by

$$\frac{1}{C(e^\varepsilon + k - 1)} \begin{cases} e^\varepsilon & \text{if } \text{HASH}_c^{(k)}(s) = y, \\ 1 & \text{if } \text{HASH}_c^{(k)}(s) \neq y. \end{cases} \quad (15)$$

For decoding, fix candidate set \mathcal{S} and interpret the privatization mechanism \mathbf{Q}_{ORR} as a $kC \times S$ row-stochastic matrix:

$$\mathbf{Q}_{\text{ORR}} = \frac{1}{C} \frac{1}{e^\varepsilon + k - 1} (\mathbf{1} + (e^\varepsilon - 1)\mathbf{H}) \quad (16)$$

where:

$$\mathbf{H}(y, c|s) = \mathbb{1}_{\{\text{HASH}_c^{(k)}(s)=y\}} \quad (17)$$

Note that \mathbf{H} is a $kC \times S$ sparse binary matrix encoding the hashed outputs for each cohort, wherein each column of \mathbf{H} has exactly C non-zero entries.

Now $\mathbf{m} = \mathbf{p}\mathbf{Q}_{\text{ORR}}$ is the expected output distribution for true probability vector \mathbf{p} , allowing us to form an empirical estimator by using standard least-squares techniques to solve the linear system:

$$\hat{\mathbf{p}}_{\text{ORR}}\mathbf{H} = \frac{1}{e^\varepsilon - 1} (C(e^\varepsilon + k - 1)\hat{\mathbf{m}} - \mathbf{1}). \quad (18)$$

Note that when $C = 1$ and \mathbf{H} is the identity matrix, (18) reduces to standard k -RR empirical estimator as seen in (6).

As with the k -RR empirical estimator, $\hat{\mathbf{p}}_{\text{ORR}}$ may have negative entries. Section 4.1 describes methods for constraining $\hat{\mathbf{p}}_{\text{ORR}}$ to \mathbb{S}^k , of which simplex projection is demonstrated to offer superior performance in Section 4.4. The remainder of the paper assumes that O-RR uses the simplex projection strategy.

5.2. O-RAPPOR

RAPPOR also extends from k -ary alphabets to open alphabets using hashing and cohorts (Erlingsson et al., 2014); we refer to this extension herein as O-RAPPOR. However, the k -RAPPOR mechanism uses a size $|\tilde{\mathcal{X}}| = 2^k$

Discrete Distribution Estimation under Local Privacy

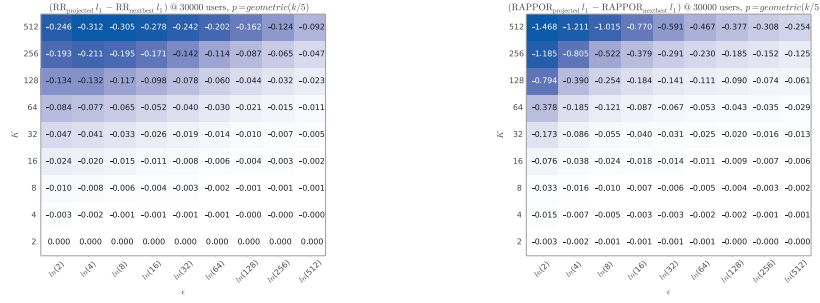


Figure 1: The improvement in ℓ_1 decoding of the projected k -RR decoder (left) and projected k -RAPPOR decoder (right). Each grid varies the size of the alphabet k (rows) and privacy parameter ϵ (columns). Each cell shows the difference in ℓ_1 magnitude that the projected decoder has over the ML and normalized k -RR decoders (left) or the standard and normalized k -RAPPOR decoders (right). Negative values mean improvement of the projected decoder over the next best alternative.

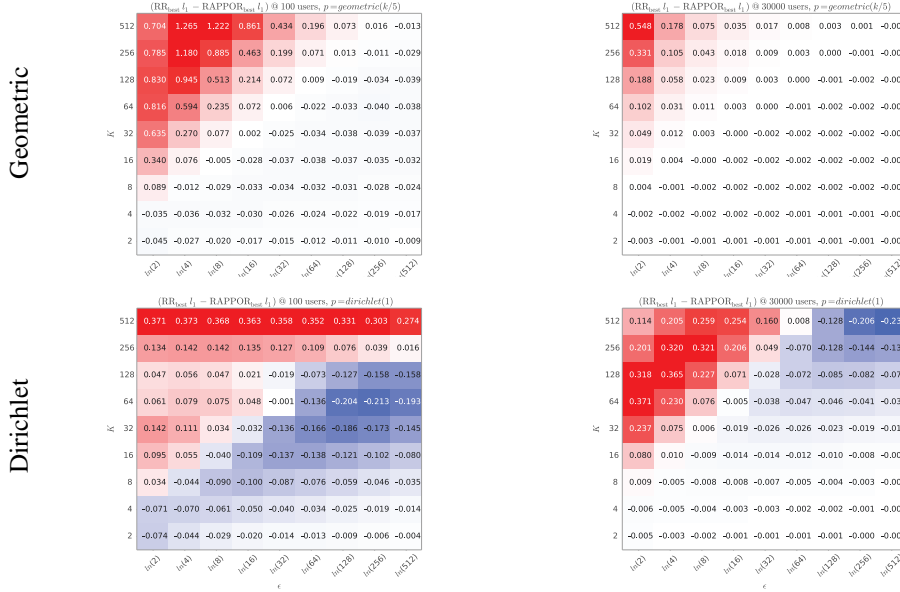


Figure 2: The improvement (negative values, blue) of the best k -RR decoder over the best k -RAPPOR decoder varying the size of the alphabet k (rows) and privacy parameter ϵ (columns). The left charts focus on small numbers of users (100); the right charts show a large number of users (30000, also representative of larger numbers of users). The top charts show the geometric distribution (skewed) and the bottom charts show the Dirichlet distribution (flat).

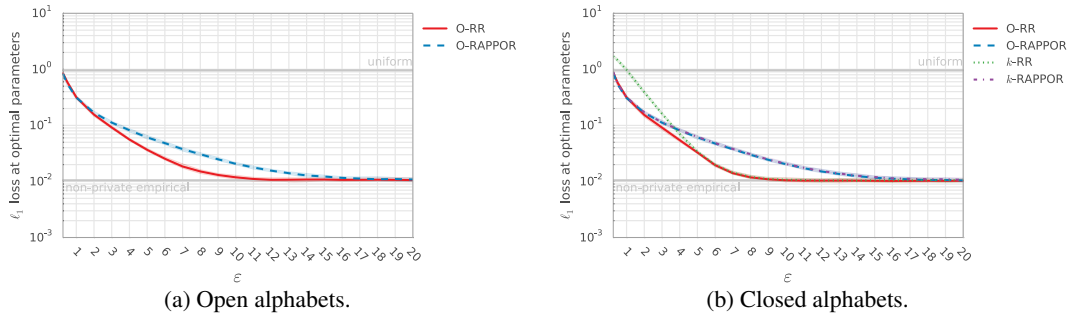


Figure 3: ℓ_1 loss of O-RR and O-RAPPOR for $n = 10^6$ on the geometric distribution when applied to unknown input alphabets (via hash functions, (a)) and to known input alphabets (via perfect hashing, (b)). Lines show median ℓ_1 loss with 90% confidence intervals over 50 samples. Free parameters are set via grid search over $k \in [2, 4, 8, \dots, 2048, 4096]$, $c \in [1, 2, 4, \dots, 512, 1024]$, $h \in [1, 2, 4, 8, 16]$ for each ϵ . Note that the k -RAPPOR and O-RAPPOR lines in (b) are nearly indistinguishable. Baselines indicate expected loss from (1) using an empirical estimator directly on the input s and (2) using the uniform distribution as the \hat{p} estimate.

input representation as opposed to k -RR’s size $|\mathcal{X}| = k$ representation. Taking advantage of the larger input space, O-RAPPOR uses an independent h -hash Bloom filter $\text{BLOOM}_c^{(k)}$ for each cohort before applying the k -RAPPOR mechanism—i.e. the j -th bit of x_i is 1 if $\text{HASH}_{c,h'}^{(k)}(s_i) = j$ for any $h' \in [1 \dots h]$, where $\text{HASH}_{c,h'}^{(k)}$ are a set of hC mutually independent hash functions modulo k .

Decoding for O-RAPPOR is described in (Erlingsson et al., 2014) and follows a similar strategy as for O-RR. However, because this paper focuses on distribution estimation rather than heavy hitter detection, we eliminate both the Lasso regression stage and filtering of imputed frequencies relative to Bonferroni corrected thresholds, retaining just the regular least-squares regression.

5.3. Simulation Analysis

We ran simulations of O-RR and O-RAPPOR for $n = 10^6$ users with input drawn from an alphabet of $S = 256$ symbols under a geometric distribution with $\text{mean} = S/5$ (see Supplementary Figure 4). As described in Section 4.4, the geometric distribution is representative of actual data and relatively easy for k -RAPPOR and challenging for k -RR. Free parameters were set to minimize the median ℓ_1 loss. Similar results for $S = 4096$ and $n = 10^6$ and 10^8 are included in the Supplementary Material.

In Figure 3(a), we see that under these conditions, O-RR matches the utility of O-RAPPOR in both the very low and high privacy regimes and exceeds the utility of O-RAPPOR over midrange privacy settings.

For O-RR, we find that the optimal k depends directly on ϵ , that increasing C consistently improves performance in the low-to-mid privacy regime, and that $C = 1$ noticeably underperforms across the range of privacy levels. For O-RAPPOR, we find that performance improves as k increases (with $k = 4096$ near the asymptotic limit), that $C = 1$ noticeably underperforms across the range of privacy values, but with all $C \geq 2$ performing indistinguishably. Finally, we find that the optimal value for h is consistently 1, indicating that Bloom filters provide no utility improvement beyond simple hashing. See Supplementary Figure 11 for details.

5.4. Improved Utility for Closed Alphabets

O-RR and O-RAPPOR extend k -ary mechanisms to open alphabets through the use of hash functions and cohorts. These same mechanisms may also be applied to closed alphabets known *a priori*. While direct application is possible, the reliance on hash functions exposes both mechanism to unnecessary risk of hash collision.

Instead, we modify the O-RR and O-RAPPOR mechanisms,

replacing each cohort’s generic hash functions with minimal perfect hash functions mapping \mathcal{S} to $[0 \dots S-1]$ before applying the modulo k operation. In most closed-alphabet applications, $\mathcal{S} = [0 \dots S-1]$, in which case these minimal perfect hash functions are simply permutations. Also note that in this setting, O-RR and O-RAPPOR reduce to exactly their k -ary counterparts when C and h are both 1 except that the output symbols are permuted.

In Figure 3(b), we evaluate these modified mechanisms using the same method described in Section 5.3 (note that the utilities of k -RAPPOR and O-RAPPOR are nearly indistinguishable). O-RAPPOR benefits little from the introduction of minimal perfect hash functions. In contrast, O-RR’s utility improves significantly, meeting or exceeding the utility of all other mechanisms at all considered ϵ .

6. Conclusion

Data improves products, services, and our understanding of the world. But its collection comes with risks to the individuals represented in the data as well as to the institutions responsible for the data’s stewardship. This paper’s focus on distribution estimation under local privacy takes one step toward a world where the benefits of data-driven insights are decoupled from the collection of raw data. Our new theoretical and empirical results show that combining cohort-style hashing with the k -ary extension of the classical randomized response mechanism admits practical, state of the art results for locally private logging.

In many applications, data is collected to enable the making of a specific decision. In such settings, the nature of the decision frequently determines the required level of utility, and the number of reports to be collected n is pre-determined by the size of the existing user base. Thus, the differential privacy practitioner’s role is often to offer users as much privacy as possible while still extracting sufficient utility at the given n . Our results suggest that O-RR may play a crucial role for such a practitioner, offering a single mechanism that provides maximal privacy at any desired utility level simply by adjusting the mechanism’s parameters.

In future work, we plan to examine estimation of non-stationary distributions as they change over time, a common scenario in data logged from user interactions. We will also consider what utility improvements may be possible when some responses need more privacy than others, another common scenario in practice. Much more work remains before we can dispel the collection of un-noised data altogether.

Acknowledgements. Thanks to Úlfar Erlingsson, Ilya Mironov, and Andrey Zhmoginov for their comments on drafts of this paper.

References

- Bassily, Raef and Smith, Adam. Local, private, efficient protocols for succinct histograms. *arXiv preprint arXiv:1504.04686*, 2015.
- Blocki, Jeremiah, Datta, Anupam, and Bonneau, Joseph. Differentially private password frequency lists. 2016.
- Boyd, Stephen and Vandenberghe, Lieven. *Convex optimization*. Cambridge university press, 2004.
- Chan, T-H Hubert, Li, Mingfei, Shi, Elaine, and Xu, Wenchang. Differentially private continual monitoring of heavy hitters from distributed streams. In *Privacy Enhancing Technologies*, pp. 140–159. Springer, 2012.
- Diakonikolas, Ilias, Hardt, Moritz, and Schmidt, Ludwig. Differentially private learning of structured discrete distributions. In *Advances in Neural Information Processing Systems*, pp. 2557–2565, 2015.
- Duchi, John, Wainwright, Martin J, and Jordan, Michael I. Local privacy and minimax bounds: Sharp rates for probability estimation. In *Advances in Neural Information Processing Systems*, pp. 1529–1537, 2013a.
- Duchi, John C, Jordan, Michael I, and Wainwright, Martin J. Local privacy, data processing inequalities, and statistical minimax rates. *arXiv preprint arXiv:1302.3203*, 2013b.
- Dwork, C. Differential privacy. In *Automata, languages and programming*, pp. 1–12. Springer, 2006.
- Dwork, C. and Lei, J. Differential privacy and robust statistics. In *Proceedings of the 41st annual ACM symposium on Theory of computing*, pp. 371–380. ACM, 2009.
- Dwork, C., McSherry, F., Nissim, K., and Smith, A. Calibrating noise to sensitivity in private data analysis. In *Theory of Cryptography*, pp. 265–284. Springer, 2006.
- Dwork, Cynthia. Differential privacy: A survey of results. In *Theory and applications of models of computation*, pp. 1–19. Springer, 2008.
- Erlingsson, Úlfar, Pihur, Vasyl, and Korolova, Aleksandra. Rappor: Randomized aggregatable privacy-preserving ordinal response. In *Proceedings of the 2014 ACM SIGSAC Conference on Computer and Communications Security*, pp. 1054–1067. ACM, 2014.
- Hsu, Justin, Khanna, Sanjeev, and Roth, Aaron. Distributed private heavy hitters. In *Automata, Languages, and Programming*, pp. 461–472. Springer, 2012.
- Kairouz, Peter, Oh, Sewoong, and Viswanath, Pramod. Extremal mechanisms for local differential privacy. In *Advances in Neural Information Processing Systems*, pp. 2879–2887, 2014.
- Kamath, Sudeep, Orlitsky, Alon, Pichapati, Venkatesh, and Suresh, Ananda Theertha. On learning distributions from their samples. In *Proceedings of The 28th Conference on Learning Theory*, pp. 1066–1100, 2015.
- Lehmann, Erich Leo and Casella, George. *Theory of point estimation*, volume 31. Springer Science & Business Media, 1998.
- McCallum, Andrew, Nigam, Kamal, et al. A comparison of event models for naive bayes text classification. In *AAAI-98 workshop on learning for text categorization*, volume 752, pp. 41–48. Citeseer, 1998.
- Narayanan, Arvind and Shmatikov, Vitaly. Robust de-anonymization of large sparse datasets. In *Security and Privacy, 2008. SP 2008. IEEE Symposium on*, pp. 111–125. IEEE, 2008.
- Wang, Weiran and Carreira-Perpiñán, Miguel Á. Projection onto the probability simplex: An efficient algorithm with a simple proof, and an application. *CoRR*, abs/1309.1541, 2013. URL <http://arxiv.org/abs/1309.1541>.
- Warner, Stanley L. Randomized response: A survey technique for eliminating evasive answer bias. *Journal of the American Statistical Association*, 60(309):63–69, 1965.

Supplementary Material: Discrete Distribution Estimation Under Local Privacy

A. Proof of Theorem 2

As argued in the proof sketch of Theorem 2, it suffices to show that $r_{\ell,\varepsilon,k,n}(\mathbf{Q})$ obeys the data processing inequality. Precisely, we need to show that for any row stochastic matrix \mathbf{W} , $r_{\ell,\varepsilon,k,n}(\mathbf{W}\mathbf{Q}) \geq r_{\ell,\varepsilon,k,n}(\mathbf{Q})$. Observe that this is equivalent to showing that $r_{\ell,\varepsilon,k,n}(\mathbf{Q}) \geq r_{\ell,k,n}$, where $r_{\ell,k,n}$ is the minimax risk in the non-private setting.

Consider the set of all randomized estimators $\hat{\mathbf{p}}$. Under randomized estimators, the minimax risk is given by

$$r_{\ell,k,n} = \inf_{\hat{\mathbf{p}}} \sup_{\mathbf{p} \in \mathbb{S}^k} \mathbb{E}_{X^n \sim \mathbf{p}, \hat{\mathbf{p}}} \ell(\mathbf{p}, \hat{\mathbf{p}}),$$

where the expectation is taken over the randomness in the observations X_1, \dots, X_n and the randomness in $\hat{\mathbf{p}}$. Under a differentially private mechanism \mathbf{Q} , the minimax risk is given by

$$r_{\ell,\varepsilon,k,n}(\mathbf{Q}) = \inf_{\hat{\mathbf{p}}_{\mathbf{Q}}} \sup_{\mathbf{p} \in \mathbb{S}^k} \mathbb{E}_{Y^n \sim \mathbf{p}\mathbf{Q}, \hat{\mathbf{p}}_{\mathbf{Q}}} \ell(\mathbf{p}, \hat{\mathbf{p}}_{\mathbf{Q}}),$$

where the expectation is taken over the randomness in the private observations Y_1, \dots, Y_n and the randomness in $\hat{\mathbf{p}}_{\mathbf{Q}}$.

Assume that there exists a (potentially randomized) estimator $\hat{\mathbf{p}}_{\mathbf{Q}}^*$ that achieves $r_{\ell,\varepsilon,k,n}(\mathbf{Q})$. Consider the following randomized estimator: \mathbf{Q} is first applied to X_1, \dots, X_n individually and $\hat{\mathbf{p}}_{\mathbf{Q}}^*$ is then jointly applied to the outputs of \mathbf{Q} . This estimator achieves a risk of $r_{\ell,\varepsilon,k,n}(\mathbf{Q})$. Therefore, $r_{\ell,k,n} \leq r_{\ell,\varepsilon,k,n}(\mathbf{Q})$.

If there is no estimator that can achieve $r_{\ell,\varepsilon,k,n}(\mathbf{Q})$, then there exists a sequence of (potentially randomized) estimators $\{\hat{\mathbf{p}}_{\mathbf{Q}}^i\}$ such that $\lim_{i \rightarrow \infty} \hat{\mathbf{p}}_{\mathbf{Q}}^i$ achieves the minimax risk. In other words, if $r_{\ell,\varepsilon,k,n}^i(\mathbf{Q})$ represents the risk under $\hat{\mathbf{p}}_{\mathbf{Q}}^i$, then $\lim_{i \rightarrow \infty} r_{\ell,\varepsilon,k,n}^i(\mathbf{Q}) = r_{\ell,\varepsilon,k,n}(\mathbf{Q})$. Using an argument similar to the one presented above, we get that $r_{\ell,k,n} \leq r_{\ell,\varepsilon,k,n}^i(\mathbf{Q})$. Taking the limit as i goes to infinity on both sides, we get that $r_{\ell,k,n} \leq r_{\ell,\varepsilon,k,n}(\mathbf{Q})$. This finishes the proof.

B. Proof of Proposition 3

Fix Q to Q_{KRR} and $\hat{\boldsymbol{p}}$ to be the empirical estimator given in (6). In this case, we have that

$$\begin{aligned}
 \mathbb{E}_{Y^n \sim \mathbf{m}(Q_{\text{KRR}})} \|\hat{\boldsymbol{p}} - \boldsymbol{p}\|_2^2 &= \mathbb{E}_{Y^n \sim \mathbf{m}(Q_{\text{KRR}})} \left\| \frac{e^\varepsilon + k - 1}{e^\varepsilon - 1} \hat{\mathbf{m}} - \frac{1}{e^\varepsilon - 1} \boldsymbol{p} \right\|_2^2 \\
 &= \mathbb{E}_{Y^n \sim \mathbf{m}(Q_{\text{KRR}})} \left\| \frac{e^\varepsilon + k - 1}{e^\varepsilon - 1} (\hat{\mathbf{m}} - \mathbf{m}) \right\|_2^2 \\
 &= \left(\frac{e^\varepsilon + k - 1}{e^\varepsilon - 1} \right)^2 \mathbb{E}_{Y^n \sim \mathbf{m}(Q_{\text{KRR}})} \|\hat{\mathbf{m}} - \mathbf{m}\|_2^2 \\
 &= \left(\frac{e^\varepsilon + k - 1}{e^\varepsilon - 1} \right)^2 \frac{1 - \sum_{i=1}^k m_i^2}{n} \\
 &= \frac{1}{n} \left(\frac{e^\varepsilon + k - 1}{e^\varepsilon - 1} \right)^2 \left(1 - \frac{\sum_{i=1}^k \{(e^\varepsilon - 1)^2 p_i^2 + 2(e^\varepsilon - 1)p_i + 1\}}{(e^\varepsilon + k - 1)^2} \right) \\
 &= \frac{(e^\varepsilon + k - 1)^2 - 2(e^\varepsilon - 1) - k - (e^\varepsilon - 1)^2 \sum_{i=1}^k p_i^2}{n(e^\varepsilon - 1)^2} \\
 &= \frac{((e^\varepsilon - 1) + k)^2 - 2(e^\varepsilon - 1) - k}{n(e^\varepsilon - 1)^2} - \frac{(e^\varepsilon - 1)^2}{n(e^\varepsilon - 1)^2} + \frac{1}{n} - \frac{\sum_{i=1}^k p_i^2}{n} \\
 &= \frac{(e^\varepsilon - 1)^2 + 2k(e^\varepsilon - 1) + k^2 - 2(e^\varepsilon - 1) - k - (e^\varepsilon - 1)^2}{n(e^\varepsilon - 1)^2} + \frac{1 - \sum_{i=1}^k p_i^2}{n} \\
 &= \frac{2(k-1)(e^\varepsilon - 1) + k(k-1)}{n(e^\varepsilon - 1)^2} + \frac{1 - \sum_{i=1}^k p_i^2}{n} \\
 &= \frac{k-1}{n} \left(\frac{2(e^\varepsilon - 1) + k}{(e^\varepsilon - 1)^2} \right) + \frac{1 - \sum_{i=1}^k p_i^2}{n},
 \end{aligned}$$

and

$$\begin{aligned}
 \mathbb{E}_{Y^n \sim \mathbf{m}(Q_{\text{KRR}})} \|\hat{\boldsymbol{p}} - \boldsymbol{p}\|_1 &= \left(\frac{e^\varepsilon + k - 1}{e^\varepsilon - 1} \right) \mathbb{E}_{Y^n \sim \mathbf{m}(Q_{\text{KRR}})} \|\hat{\mathbf{m}} - \mathbf{m}\|_1 \\
 &= \left(\frac{e^\varepsilon + k - 1}{e^\varepsilon - 1} \right) \sum_{i=1}^k \mathbb{E} |m_i - \hat{m}_i| \\
 &\approx \left(\frac{e^\varepsilon + k - 1}{e^\varepsilon - 1} \right) \sum_{i=1}^k \sqrt{\frac{2m_i(1-m_i)}{\pi n}} \\
 &= \frac{1}{e^\varepsilon - 1} \sum_{i=1}^k \sqrt{\frac{2((e^\varepsilon - 1)p_i + 1)((e^\varepsilon - 1)(1-p_i) + k - 1)}{\pi n}}.
 \end{aligned}$$

C. Proof of Proposition 4

Fix Q to $Q_{k\text{-RAPPOR}}$ and \hat{p} to be the empirical estimator given in (11), and let $C = \frac{e^{\varepsilon/2}-1}{e^{\varepsilon/2}+1}$, $B = \frac{1}{e^{\varepsilon/2}+1}$, and $A = e^{\varepsilon/2} - 1$. Then $C = BA$, $1 - B = e^{\varepsilon/2}B$, and from Section 4.2 $m_i = p_i C + B$. Using this notation, we have that

$$\begin{aligned}
 \mathbb{E}_{Y^n \sim \mathbf{m}(Q_{k\text{-RAPPOR}})} \|\hat{p} - \mathbf{p}\|_2^2 &= \mathbb{E}_{Y^n \sim \mathbf{m}(Q_{k\text{-RAPPOR}})} \left\| \frac{e^{\varepsilon/2} + 1}{e^{\varepsilon/2} - 1} \hat{\mathbf{m}} - \frac{1}{e^{\varepsilon/2} - 1} \mathbf{p} \right\|_2^2 \\
 &= \mathbb{E}_{Y^n \sim \mathbf{m}(Q_{k\text{-RAPPOR}})} \left\| \frac{e^{\varepsilon/2} + 1}{e^{\varepsilon/2} - 1} (\hat{\mathbf{m}} - \mathbf{m}) \right\|_2^2 \\
 &= \left(\frac{e^{\varepsilon/2} + 1}{e^{\varepsilon/2} - 1} \right)^2 \mathbb{E}_{Y^n \sim \mathbf{m}(Q_{k\text{-RAPPOR}})} \|\hat{\mathbf{m}} - \mathbf{m}\|_2^2 \\
 &= \frac{1}{nC^2} \left(C + kB - \sum_{i=1}^k (p_i C + B)^2 \right) \\
 &= \frac{1}{n} \left(1 - \sum_{i=1}^k p_i^2 \right) + \frac{1}{nC^2} (C - C^2 + kB - kB^2 - 2CB) \\
 &= \frac{1}{n} \left(1 - \sum_{i=1}^k p_i^2 \right) + \frac{1}{nBA^2} (A - BA^2 + k(1 - B) - 2BA) \\
 &= \frac{1}{n} \left(1 - \sum_{i=1}^k p_i^2 \right) + \frac{1}{n} \frac{ke^{\varepsilon/2}}{(e^{\varepsilon/2} - 1)^2},
 \end{aligned}$$

and

$$\begin{aligned}
 \mathbb{E}_{Y^n \sim \mathbf{m}(Q_{\text{KRR}})} \|\hat{p} - \mathbf{p}\|_1 &= \left(\frac{e^{\varepsilon/2} + 1}{e^{\varepsilon/2} - 1} \right) \mathbb{E}_{Y^n \sim \mathbf{m}(Q_{k\text{-RAPPOR}})} \|\hat{\mathbf{m}} - \mathbf{m}\|_1 \\
 &= \left(\frac{e^{\varepsilon/2} + 1}{e^{\varepsilon/2} - 1} \right) \sum_{i=1}^k \mathbb{E} |m_i - \hat{m}_i| \\
 &\approx \left(\frac{e^{\varepsilon/2} + 1}{e^{\varepsilon/2} - 1} \right) \sum_{i=1}^k \sqrt{\frac{2m_i(1 - m_i)}{\pi n}} \\
 &= \sum_{i=1}^k \sqrt{\frac{2((e^{\varepsilon/2} - 1)p_i + 1)((e^{\varepsilon/2} - 1)(1 - p_i) + 1)}{\pi n(e^{\varepsilon/2} - 1)^2}}.
 \end{aligned}$$

D. Proof of Proposition 5

We want to show that for all $\mathbf{p} \in \mathbb{S}^k$ and all $\varepsilon \geq \ln k$,

$$\mathbb{E} \|\hat{p}_{\text{KRR}} - \mathbf{p}\|_2^2 \leq \mathbb{E} \|\hat{p}_{\text{RAPPOR}} - \mathbf{p}\|_2^2, \quad (19)$$

where \hat{p}_{KRR} is the empirical estimate of \mathbf{p} under $k\text{-RR}$, \hat{p}_{RAPPOR} is the empirical estimate of \mathbf{p} under $k\text{-RAPPOR}$, and \hat{p} is the empirical estimator under $k\text{-RAPPOR}$.

From propositions 3 and 4, we have that

$$\mathbb{E} \|\hat{p}_{\text{KRR}} - \mathbf{p}\|_2^2 = \frac{1 - \sum_{i=1}^k p_i^2}{n} + \frac{k-1}{n} \left(\frac{2}{e^\varepsilon - 1} + \frac{k}{(e^\varepsilon - 1)^2} \right),$$

and

$$\mathbb{E} \|\hat{p}_{\text{RAPPOR}} - \mathbf{p}\|_2^2 = \frac{1 - \sum_{i=1}^k p_i^2}{n} + \frac{ke^{\varepsilon/2}}{n(e^{\varepsilon/2} - 1)^2}.$$

Therefore, we just have to prove that

$$(k-1) \left(\frac{2}{e^\varepsilon - 1} + \frac{k}{(e^\varepsilon - 1)^2} \right) \leq \frac{ke^{\varepsilon/2}}{(e^{\varepsilon/2} - 1)^2},$$

for $\varepsilon \geq \ln k$. Alternatively, we can show that

$$f(\varepsilon, k) = \frac{k}{k-1} \left(\frac{e^\varepsilon - 1}{e^{\varepsilon/2} - 1} \right)^2 \frac{e^{\varepsilon/2}}{2e^\varepsilon + k - 2} \geq 1,$$

for $\varepsilon \geq \ln k$. Observe that $f(\varepsilon, k)$ is an increasing function of ε and therefore, it suffices to show that

$$f(\ln k, k) = \frac{k}{k-1} \left(\frac{k-1}{\sqrt{k}-1} \right)^2 \frac{\sqrt{k}}{3k-2} = \frac{k}{3k-2} \frac{\sqrt{k}(k-1)}{(\sqrt{k}-1)^2} \geq 1. \quad (20)$$

As a discrete function of $k \in \{2, 3, \dots\}$, $f(\ln k, k)$ admits a unique minimum at $k = 7$. Therefore, we just need to verify that $f(\ln 7, 7) > 1$. Indeed, $f(\ln 7, 7) = 3.1559 > 1$.

E. Discrete Distribution Estimation

Consider the $(k-1)$ -dimensional probability simplex

$$\mathbb{S}^k = \{\mathbf{p} = (p_1, \dots, p_k) \mid p_i \geq 0, \sum_{i=1}^k p_i = 1\}.$$

The discrete distribution estimation problem is defined as follows. Given a vector $\mathbf{p} \in \mathbb{S}^k$, samples X_1, \dots, X_n are drawn i.i.d according to \mathbf{p} . Our goal is to estimate the probability vector \mathbf{p} from the observation vector $X^n = (X_1, \dots, X_n)$.

An estimator $\hat{\mathbf{p}}$ is a mapping from X^n to a point in \mathbb{S}^k . The performance of $\hat{\mathbf{p}}$ may be measured via a loss function ℓ that computes a distance-like metric between $\hat{\mathbf{p}}$ and \mathbf{p} . Common loss functions include, among others, the absolute error loss $\ell_1(\mathbf{p}, \hat{\mathbf{p}}) = \sum_{i=1}^k |p_i - \hat{p}_i|$ and the quadratic loss $\ell_2^2(\mathbf{p}, \hat{\mathbf{p}}) = \sum_{i=1}^k (p_i - \hat{p}_i)^2$. The choice of the loss function depends on the application; for example, ℓ_1 loss is commonly used in classification and other machine learning applications. Given a loss function ℓ , the expected loss under $\hat{\mathbf{p}}$ after observing n i.i.d samples is given by

$$r_{\ell, k, n}(\mathbf{p}, \hat{\mathbf{p}}) = \mathbb{E}_{X^n \sim \text{Multinomial}(n, \mathbf{p})} \ell(\mathbf{p}, \hat{\mathbf{p}}). \quad (21)$$

E.1. Maximum likelihood and empirical estimation

In the absence of a prior on \mathbf{p} , a natural and commonly used estimator of \mathbf{p} is the maximum likelihood (ML) estimator. The maximum likelihood estimate $\hat{\mathbf{p}}_{\text{ML}}$ of \mathbf{p} is defined as

$$\hat{\mathbf{p}}_{\text{ML}} = \operatorname{argmax}_{\mathbf{p} \in \mathbb{S}^k} \mathbb{P}(X_1, \dots, X_n \mid \mathbf{p})$$

In this setting, it is easy to show that the maximum likelihood estimate is equivalent to the empirical estimator of \mathbf{p} , given by $\hat{p}_i = T_i/n$ where T_i is the frequency of element i . Observe that the empirical estimator is an unbiased estimator for \mathbf{p} because $\mathbb{E}[\hat{p}_i] = p_i$ for any k, n , and i . Under maximum likelihood estimation, the ℓ_2^2 loss is the most tractable and simplest to analyze loss function. Because $T_i \sim \text{Binomial}(p_i, n)$, we have $\mathbb{E}[T_i] = np_i$, $\text{Var}(T_i) = np_i(1 - p_i)$, and the expected ℓ_2^2 loss of the empirical estimator is given by

$$\begin{aligned} r_{\ell_2^2, k, n}(\mathbf{p}, \hat{\mathbf{p}}_{\text{ML}}) &= \mathbb{E} \|\hat{\mathbf{p}}_{\text{ML}} - \mathbf{p}\|_2^2 = \sum_{i=1}^k \mathbb{E} \left(\frac{T_i}{n} - p_i \right)^2 \\ &= \sum_{i=1}^k \frac{\text{Var}(T_i)}{n^2} = \frac{1 - \sum_{i=1}^k p_i^2}{n}. \end{aligned}$$

Let $\mathbf{p}_U = (\frac{1}{k}, \dots, \frac{1}{k})$ and observe that

$$r_{\ell_2^2, k, n}(\mathbf{p}, \hat{\mathbf{p}}_{\text{ML}}) \leq r_{\ell_2^2, k, n}(\mathbf{p}_U, \hat{\mathbf{p}}_{\text{ML}}) = \frac{1 - \frac{1}{k}}{n}. \quad (22)$$

In other words, the uniform distribution is the worst distribution for the empirical estimator under the ℓ_2^2 loss. From (Kamath et al., 2015), the asymptotic performance of the empirical estimator under the ℓ_1 loss functions is given by

$$r_{\ell_1, k, n}(\mathbf{p}, \hat{\mathbf{p}}_{\text{ML}}) \approx \sum_{i=1}^k \sqrt{\frac{2p_i(1-p_i)}{\pi n}},$$

where $a_n \approx b_n$ means $\lim_{n \rightarrow \infty} a_n/b_n = 1$. As in the ℓ_2^2 case, notice that

$$r_{\ell_1, k, n}(\mathbf{p}, \hat{\mathbf{p}}_{\text{ML}}) \leq r_{\ell_1, k, n}^{\ell_1}(\mathbf{p}_U, \hat{\mathbf{p}}_{\text{ML}}) \approx \sqrt{\frac{2(k-1)}{\pi n}}, \quad (23)$$

for any $\mathbf{p} \in \mathbb{S}^k$. In other words, the uniform distribution is the worst distribution for the empirical estimator under the ℓ_1 loss as well. Observe that the ℓ_1 loss scales as $\sqrt{k/n}$ whereas the ℓ_2^2 loss scales as $1/n$.

E.2. Minimax estimation

Another popular estimator that is widely studied in the absence of a prior is the minimax estimator $\hat{\mathbf{p}}_{\text{MM}}$. The minimax estimator minimizes the expected loss under the worst distribution \mathbf{p} :

$$\hat{\mathbf{p}}_{\text{MM}} = \operatorname{argmin}_{\hat{\mathbf{p}}} \max_{\mathbf{p} \in \mathbb{S}^k} \mathbb{E}_{X^n \sim \mathbf{p}} \ell(\mathbf{p}, \hat{\mathbf{p}}). \quad (24)$$

The minimax risk is therefore defined as

$$r_{\ell, k, n} = \min_{\hat{\mathbf{p}}} \max_{\mathbf{p} \in \mathbb{S}^k} \mathbb{E}_{X^n \sim \mathbf{p}} \ell(\mathbf{p}, \hat{\mathbf{p}}).$$

For the ℓ_2^2 loss, it is shown in (Lehmann & Casella, 1998) that

$$\hat{p}_i = \frac{\frac{\sqrt{n}}{k} + \sum_{j=1}^n \mathbb{1}_{\{X_j=i\}}}{\sqrt{n} + n} = \frac{\frac{\sqrt{n}}{k} + T_i}{\sqrt{n} + n}, \quad (25)$$

is the minimax estimator, and that the minimax risk is

$$r_{\ell_2^2, k, n} = \frac{1 - \frac{1}{k}}{(\sqrt{n} + 1)^2}. \quad (26)$$

Observe that unlike the empirical estimator, the minimax estimator is not even asymptotically unbiased. Moreover, it improves on the empirical estimator only slightly (compare Equations (22) to (26)), increasing the denominator from n to $n + 2\sqrt{n} + 1$ under the worst case distribution (the uniform distribution). This explains why the minimax estimator is almost never used in practice.

The minimax estimator under ℓ_1 loss is not known. However, the minimax risk is known for the case when k is fixed and n is increased. In this case, it is shown in (Kamath et al., 2015) that

$$r_{\ell_1, k, n} = \sqrt{\frac{2(k-1)}{\pi n}} + O\left(\frac{1}{n^{3/4}}\right). \quad (27)$$

Comparing Equations (23) to (27), we see that the worst case loss under the empirical estimator is again roughly as good as the minimax risk.

F. Maximum Likelihood Estimation for k -ary Mechanisms

F.1. k -RR

Proposition 6 *The maximum likelihood estimator of \mathbf{p} under k -RR is given by*

$$\hat{p}_i = \left[\frac{T_i}{\lambda} - \frac{1}{e^\varepsilon - 1} \right]^+, \quad (28)$$

where $[x]^+ = \max(0, x)$, T_i is the frequency of element i calculated from Y^n , and λ is chosen so that

$$\sum_{i=1}^k \left[\frac{T_i}{\lambda} - \frac{1}{e^\varepsilon - 1} \right]^+ = 1. \quad (29)$$

Moreover, finding λ can be done in $O(k \log k)$ steps.

The proof of the above proposition is provided in Supplementary Section F.2.

F.2. Proof of Proposition 6

The maximum likelihood estimator under k -RR is the solution to

$$\hat{\mathbf{p}}_{\text{ML}} = \operatorname{argmax}_{\mathbf{p} \in \mathbb{S}^k} \mathbb{P}(Y_1, \dots, Y_n | \mathbf{p}),$$

where the Y_i 's are the outputs of k -RR. Since the $\log(\cdot)$ function is a monotonic function, the above maximum likelihood estimation problem is equivalent to

$$\hat{\mathbf{p}}_{\text{ML}} = \operatorname{argmax}_{\mathbf{p} \in \mathbb{S}^k} \log \mathbb{P}(Y_1, \dots, Y_n | \mathbf{p}).$$

Given that

$$\begin{aligned} \mathbb{P}(Y_1, \dots, Y_n | \mathbf{p}) &= \prod_{i=1}^n \mathbb{P}(Y_i | \mathbf{p}) \\ &= \prod_{i=1}^n \left(\sum_{j=1}^k \mathbf{Q}_{\text{KRR}}(Y_i | X_i = j) p_j \right), \end{aligned}$$

we have that

$$\log \mathbb{P}(Y_1, \dots, Y_n | \mathbf{p}) = \sum_{i=1}^n \log \left(\sum_{j=1}^k \mathbf{Q}_{\text{KRR}}(Y_i | X_i = j) p_j \right).$$

Observe that

$$\sum_{j=1}^k \mathbf{Q}_{\text{KRR}}(Y_i | X_i = j) p_j = \mathbf{Q}_{\text{KRR}}(Y_i | X_i = Y_i) p_{Y_i} + \sum_{j \neq Y_i} \mathbf{Q}_{\text{KRR}}(Y_i | X_i = j) p_j \quad (30)$$

$$= \frac{e^\varepsilon}{e^\varepsilon + k - 1} p_{Y_i} + \frac{1}{e^\varepsilon + k - 1} (1 - p_{Y_i}) \quad (31)$$

$$= \frac{1}{e^\varepsilon + k - 1} ((e^\varepsilon - 1) p_{Y_i} + 1), \quad (32)$$

and therefore,

$$\sum_{i=1}^n \log \left(\sum_{j=1}^k \mathbf{Q}_{\text{KRR}}(Y_i | X_i = j) p_j \right) = \sum_{i=1}^n T_i \log \left(\frac{1}{e^\varepsilon + k - 1} ((e^\varepsilon - 1) p_i + 1) \right),$$

where T_i is the number of Y 's that are equal to i (i.e., the frequency of element i in the observed sequence Y^n). Thus, the maximum likelihood estimation problem under k -RR is equivalent to

$$\hat{\mathbf{p}}_{\text{ML}} = \operatorname{argmax}_{\mathbf{p} \in \mathbb{S}^k} \sum_{i=1}^k T_i \log((e^\varepsilon - 1)p_i + 1).$$

The above constrained optimization problem is a convex optimization problem that is well studied in the literature under the rubric of water-filling algorithms. From (Boyd & Vandenberghe, 2004), the solution to this problem is given by

$$\hat{p}_i = \left[\frac{T_i}{\lambda} - \frac{1}{e^\varepsilon - 1} \right]^+,$$

where $[x]^+ = \max(0, x)$ and λ is chosen so that

$$\sum_{i=1}^k \left[\frac{T_i}{\lambda} - \frac{1}{e^\varepsilon - 1} \right]^+ = 1.$$

Given the T_i 's, \mathbf{p} is computed according to the empirical estimator. If all the \hat{p}_i 's are non-negative, then the maximum likelihood estimate is the same as the empirical estimate. If not, $\hat{\mathbf{p}}$ is sorted, its negative entries are zeroed out, and lambda is computed according to the above equation. Given lambda, a new $\hat{\mathbf{p}}$ can be computed and the above process can be repeated until all the entries of $\hat{\mathbf{p}}$ are non-negative. Notice that sorting happens once and the process is repeated at most $k-1$ times. Therefore, the computational complexity of this algorithm is upper bounded by $k \log k + k$ which is $O(k \log k)$.

F.3. k -RAPPOR

Proposition 7 *The maximum likelihood estimator of \mathbf{p} under k -RAPPOR is*

$$\operatorname{argmax}_{\mathbf{p} \in \mathbb{S}^k} \sum_{j=1}^k (n - T_j) \log((1 - \delta) - (1 - 2\delta)p_j) + T_j \log((1 - 2\delta)p_j + \delta)$$

where $T_j = \sum_{i=1}^n Y_i^{(j)}$ and $\delta = 1/(e^{\varepsilon/2} + 1)$.

The proof of the above proposition is provided in Supplementary Section F.4. Observe that unlike k -RR, a k -dimensional convex program has to be solved in this case to determine the maximum likelihood estimate of \mathbf{p} .

F.4. Proof of Proposition 7

The maximum likelihood estimator under k -RAPPOR is the solution to

$$\hat{\mathbf{p}}_{\text{ML}} = \operatorname{argmax}_{\mathbf{p} \in \mathbb{S}^k} \mathbb{P}(Y_1, \dots, Y_n | \mathbf{p}),$$

where the Y_i 's are the outputs of k -RAPPOR. Since the $\log(\cdot)$ function is a monotonic function, the above maximum likelihood estimation problem is equivalent to

$$\hat{\mathbf{p}}_{\text{ML}} = \operatorname{argmax}_{\mathbf{p} \in \mathbb{S}^k} \log \mathbb{P}(Y_1, \dots, Y_n | \mathbf{p}).$$

Recall that under k -RAPPOR, $Y_i = [Y_i^{(1)}, \dots, Y_i^{(k)}]$ is a k -dimensional binary vector, which implies that

$$\mathbb{P}(Y_i^{(j)} = 1) = \left(\frac{e^{\varepsilon/2} - 1}{e^{\varepsilon/2} + 1} \right) p_j + \frac{1}{e^{\varepsilon/2} + 1}, \quad (33)$$

for all $i \in \{1, \dots, n\}$ and $j \in \{1, \dots, k\}$. Therefore,

$$\begin{aligned} \log \mathbb{P}(Y_1, \dots, Y_n | \mathbf{p}) &= \log \prod_{i=1}^n \prod_{j=1}^k \left(Y_i^{(j)} (p_j(1-\delta) + (1-p_j)\delta) + (1-Y_i^{(j)}) (p_j\delta + (1-p_j)(1-\delta)) \right) \\ &= \sum_{i=1}^n \sum_{j=1}^k \log \left(Y_i^{(j)} (p_j(1-\delta) + (1-p_j)\delta) + (1-Y_i^{(j)}) (p_j\delta + (1-p_j)(1-\delta)) \right) \\ &= \sum_{i=1}^n \sum_{j=1}^k \log \left((1-2\delta)(2Y_i^{(j)} - 1)p_j - Y_i^{(j)}(1-2\delta) + (1-\delta) \right), \end{aligned}$$

where $\delta = 1/(1 + e^{\varepsilon/2})$. Therefore, under k -RAPPOR, the maximum likelihood estimation problem is given by

$$\operatorname{argmax}_{\mathbf{p} \in \mathbb{S}^k} \sum_{j=1}^k (n - T_j) \log((1-\delta) - (1-2\delta)p_j) + T_j \log((1-2\delta)p_j + \delta)$$

where $T_j = \sum_{i=1}^n Y_i^{(j)}$.

G. Conditions for Accurate Decoding under k -RR

For accurate decoding, we must satisfy three criteria: (i) k and C must be large enough that the input strings to be distinguishable, (ii) k and C must be large enough that the linear system in (18) is not underconstrained, and (iii) n must be large enough that the variance on estimated probability vector $\hat{\mathbf{p}}$ is small.

Let us first consider string distinguishability. Each string $s \in \mathcal{S}$ is associated with a C -tuple of hashes it can produce in the various cohorts: $\text{HASH}_C^{(k)}(s) = \langle \text{HASH}_1^{(k)}(s), \text{HASH}_2^{(k)}(s), \dots, \text{HASH}_C^{(k)}(s) \rangle \in \mathcal{X}^C$. Two strings $s_i \in \mathcal{S}$ and $s_j \in \mathcal{S}$ are distinguishable from one another under the encoding scheme if $\text{HASH}_C^{(k)}(s_i) \neq \text{HASH}_C^{(k)}(s_j)$, and a string s is distinguishable within the set \mathcal{S} if $\text{HASH}_C^{(k)}(s) \neq \text{HASH}_C^{(k)}(s_j) \forall s_j \in \mathcal{S} \setminus s$.

Because $\text{HASH}_C^{(k)}(s)$ is distributed uniformly over \mathcal{X}^C , $\mathbb{P}(\text{HASH}_C^{(k)}(s) = \mathbf{x}_C) \approx \frac{1}{k^C}$ for all $\mathbf{x}_C \in \mathcal{X}^C$. It follows that the probability of two strings being distinguishable is also $\frac{1}{k^C}$. Furthermore, the probability that exactly one string from \mathcal{S} produces the hash tuple \mathbf{x}_C is:

$$\text{Binomial}(1; \frac{1}{k^C}, S) = \frac{S(k^C - 1)^{S-1}}{(k^C)^S}$$

Thus, the expected number of $\mathbf{x}_C \in \mathcal{X}^C$ associated with exactly one string in \mathcal{S} , which is also the expected number of distinguishable strings in a set \mathcal{S} is:

$$\sum_{\mathbf{x}_C \in \mathcal{Y}^C} \left(\frac{S(k^C - 1)^{S-1}}{(k^C)^S} \right) = S \left(\frac{k^C - 1}{k^C} \right)^{S-1} \quad (34)$$

and the probability that a string s is distinguishable within the set \mathcal{S} is $\left(\frac{k^C - 1}{k^C} \right)^{S-1}$.

Consider a probability distribution $\mathbf{p} \in \mathbb{S}^S$. The expected recoverable probability mass is the the mass associated with the distinguishable strings within the set \mathcal{S} is $\sum_{s \in \mathcal{S}} p_s \left(\frac{k^C - 1}{k^C} \right)^{S-1} = \left(\frac{k^C - 1}{k^C} \right)^{S-1}$. Therefore, if we hope to recover at least P_t of the probability mass, we require $\left(\frac{k^C - 1}{k^C} \right)^{S-1} \geq P_t$, or equivalently, $k^C \geq \frac{1}{1 - P_t^{\frac{1}{S-1}}}$.

Now consider ensuring that the linear system in (18) is not underconstrained. The system has S variables and kC independent equations. Thus, the system is not underconstrained so long as $kC \geq S$.

H. Supplementary Figures

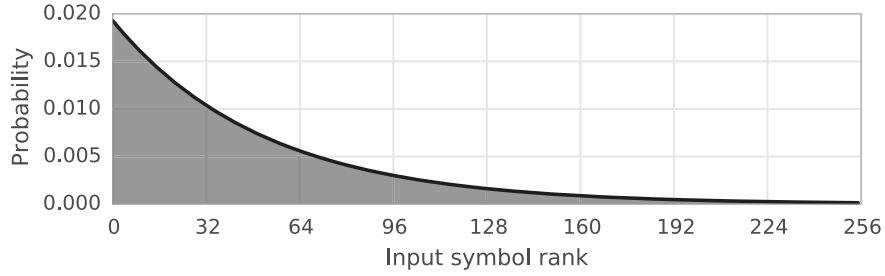


Figure 4: The true input distribution p for open-set and closed-set experiments in sections 4.4 and 5 is the geometric distribution with mean at $|\text{input alphabet}|/5$, truncated and renormalized. In the k -ary experiments of Section 4.4, the input alphabet is size k ; in the open alphabet experiments of Section 5, the input alphabet is size $S = 256$.

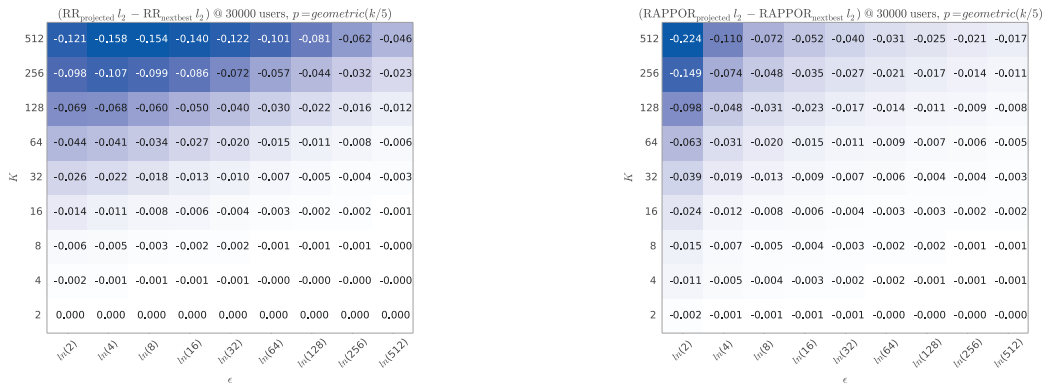


Figure 5: The improvement in l_2 decoding of the projected k -RR decoder (left) and projected k -RAPPOR decoder (right). This figure demonstrates that the same patterns hold in l_2 as in l_1 for the conditions shown in Figure 1.

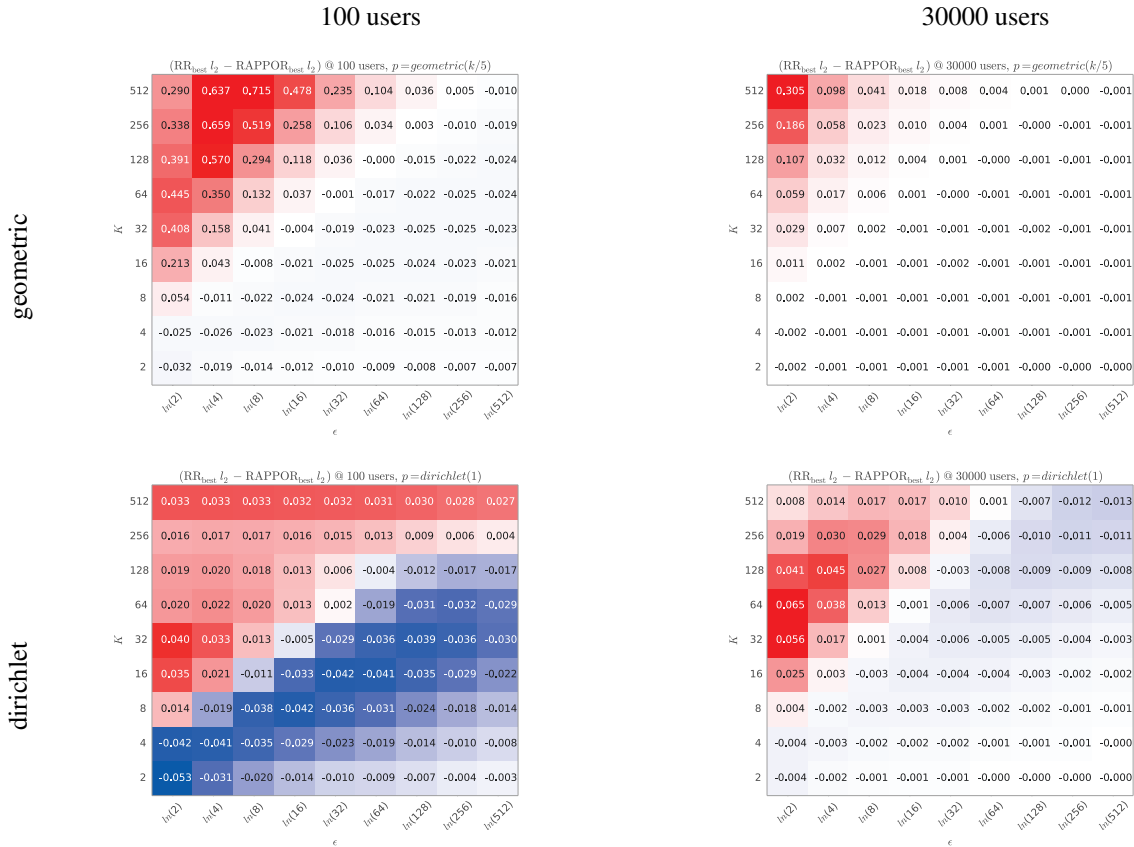


Figure 6: The improvement (negative values, blue) of the best k -RR decoder over the best k -RAPPOR decoder varying the size of the alphabet k (rows) and privacy parameter ϵ (columns). This figure demonstrates that the same patterns hold in ℓ_2 as in ℓ_1 for the conditions shown in Figure 2.

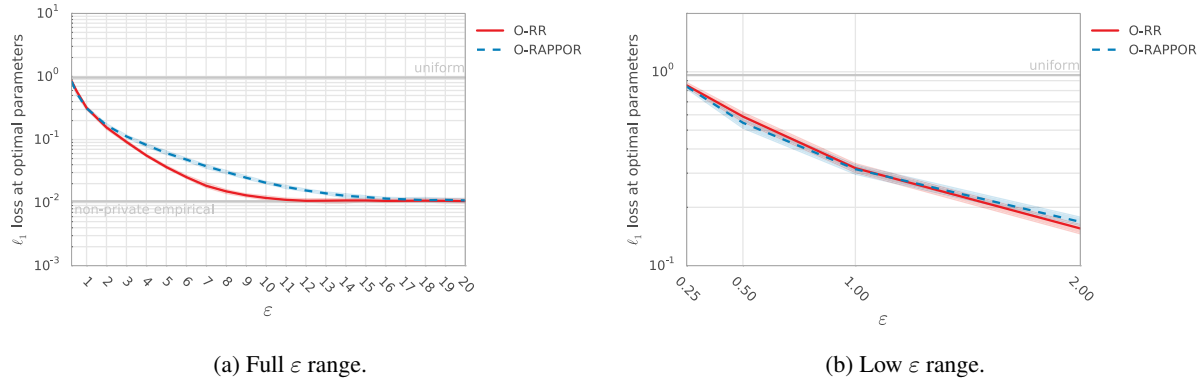


Figure 7: ℓ_1 loss when decoding open alphabets using the O-RR and O-RAPPOR for $n = 10^6$ users with input drawn from an alphabet of $S = 256$ symbols under a geometric distribution with mean= $S/5$, as depicted in Figure 4. Free parameters are set via grid search over $k \in [2, 4, 8, \dots, 2048, 4096]$, $c \in [1, 2, 4, \dots, 512, 1024]$, $h \in [1, 2, 4, 8, 16]$ to minimize the median loss over 50 samples at the given ϵ value. Lines show median ℓ_1 loss while the (narrow) shaded regions indicate 90% confidence intervals (over 50 samples). Baselines indicate expected loss from (1) using an empirical estimator directly on the input s and (2) using the uniform distribution as the \hat{p} estimate.

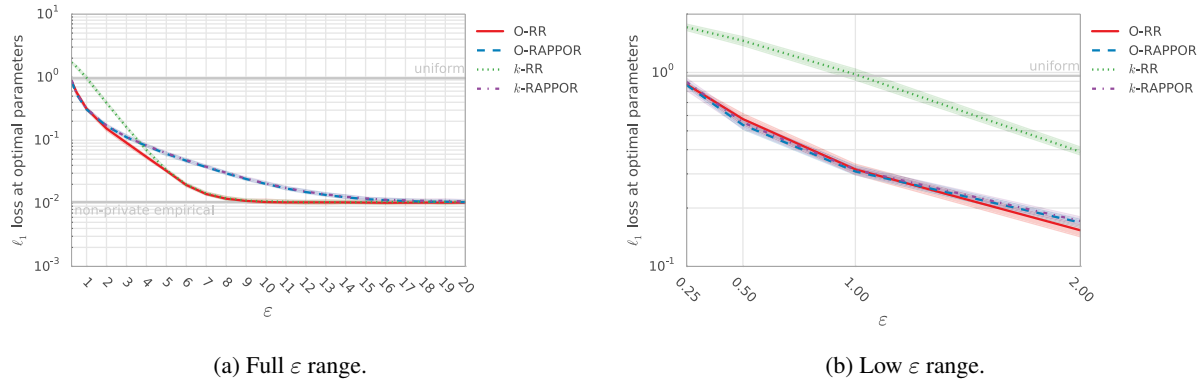


Figure 8: ℓ_1 loss when decoding a known alphabet using the O-RR and O-RAPPOR (via permutative perfect hash functions) for $n = 10^6$ users with input drawn from an alphabet of $S = 256$ symbols under a geometric distribution with mean= $S/5$, as depicted in Figure 4. Free parameters are set via grid search over $k \in [2, 4, 8, \dots, 2048, 4096]$, $c \in [1, 2, 4, \dots, 512, 1024]$, $h \in [1, 2, 4, 8, 16]$ to minimize the median loss over 50 samples at the given ϵ value. Lines show median ℓ_1 loss while the (narrow) shaded regions indicate 90% confidence intervals (over 50 samples). Note that the k -RAPPOR and O-RAPPOR lines in (b) are nearly indistinguishable. Baselines indicate expected loss from (1) using an empirical estimator directly on the input s and (2) using the uniform distribution as the \hat{p} estimate.

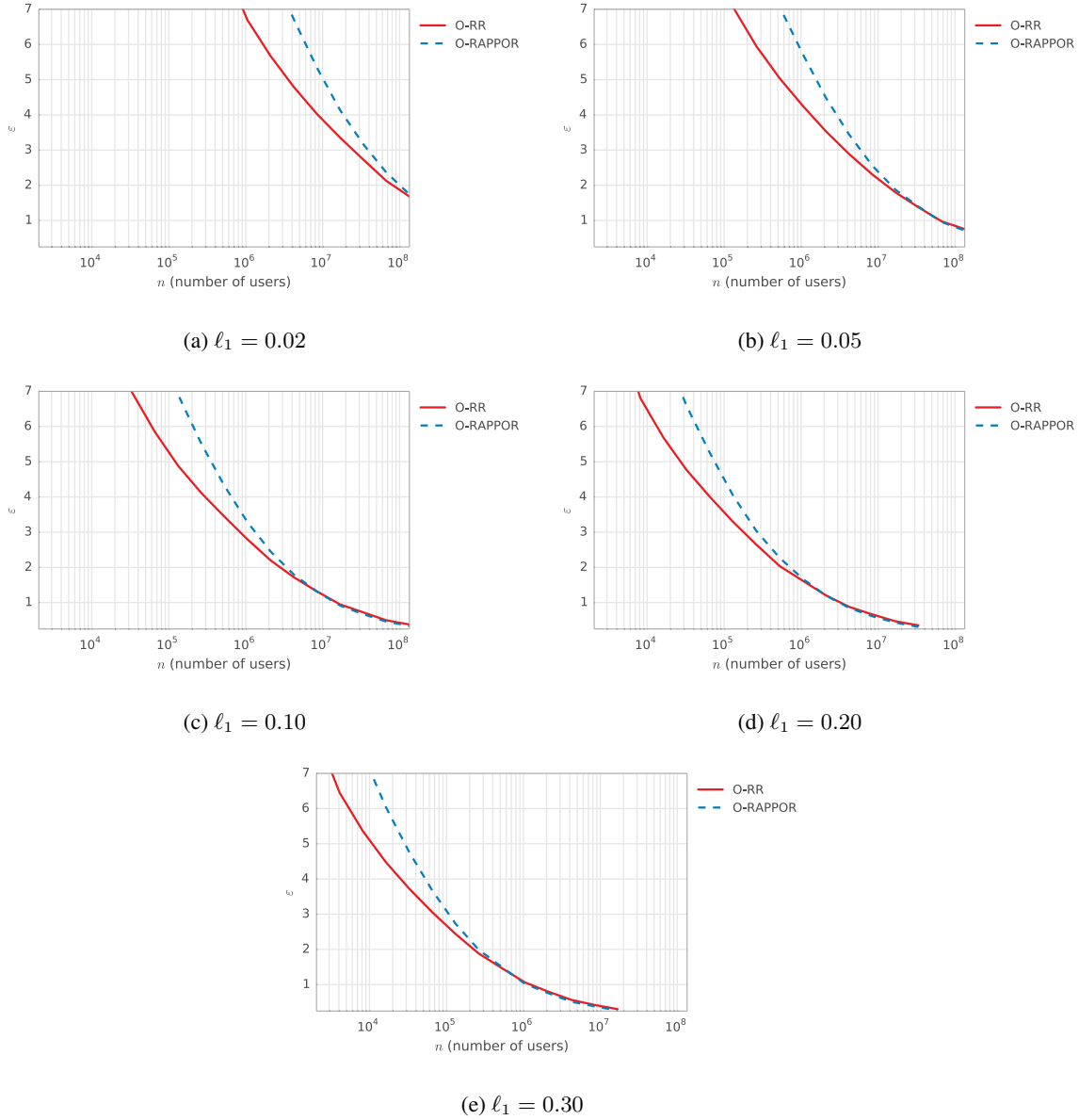


Figure 9: Taking ℓ_1 loss (the utility) and n (the number of users) as fixed requirements (as is the case in many practical scenarios), we approximate the degree of privacy ε that can be obtained under O-RR and O-RAPPOR for open alphabets (lower ε is better). Input is generated from an alphabet of $S = 256$ symbols under a geometric distribution with mean= $S/5$, as depicted in Figure 4. Free parameters are set via grid search to minimize the median loss over 50 samples at the given ε and fixed parameter values.

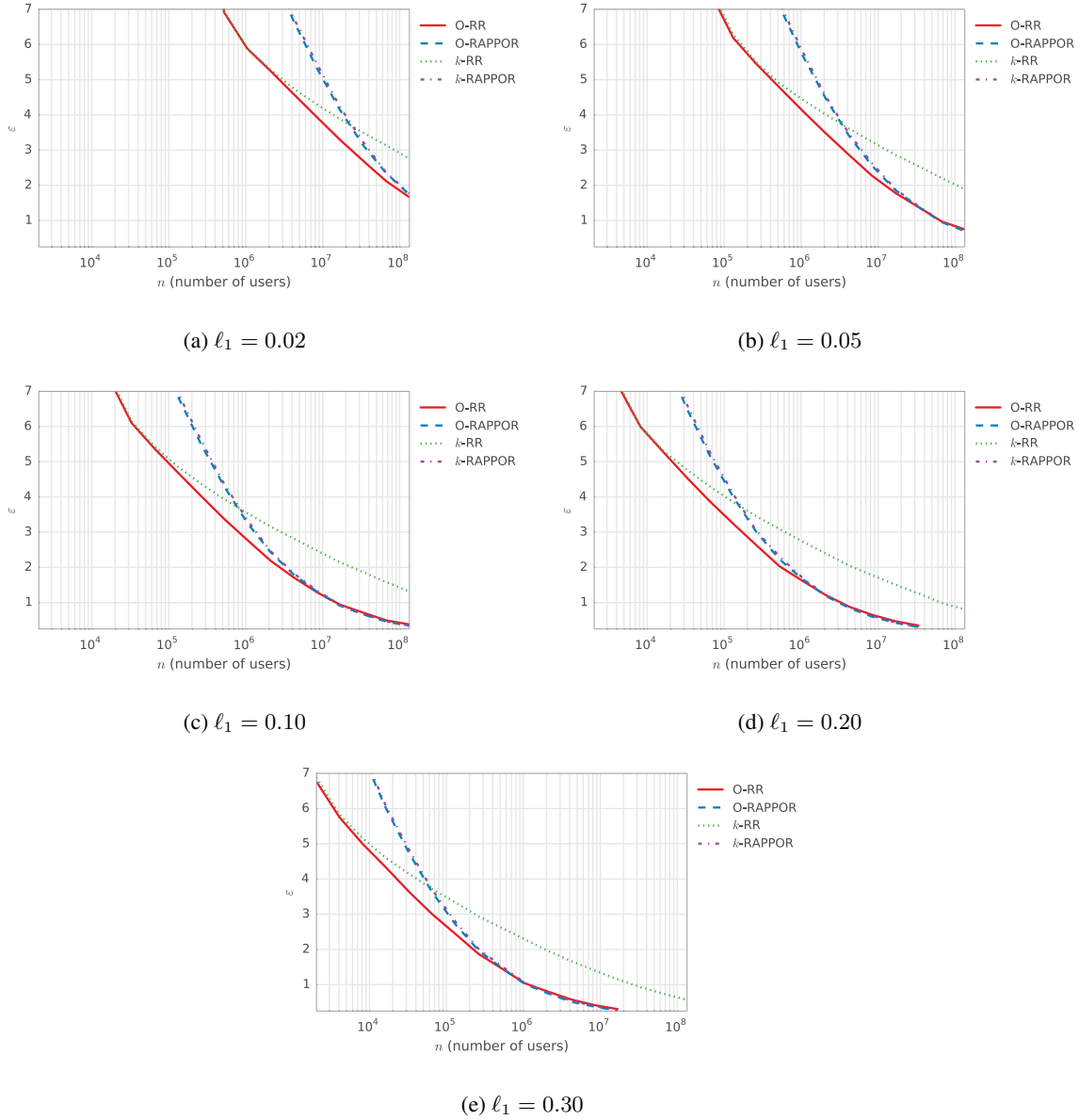
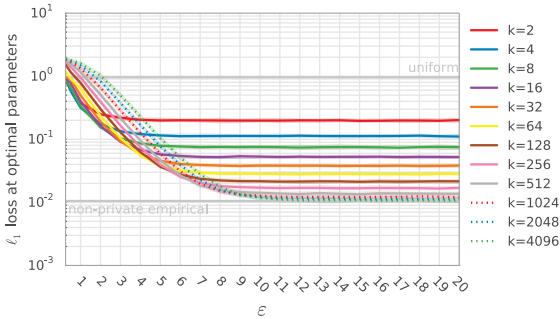
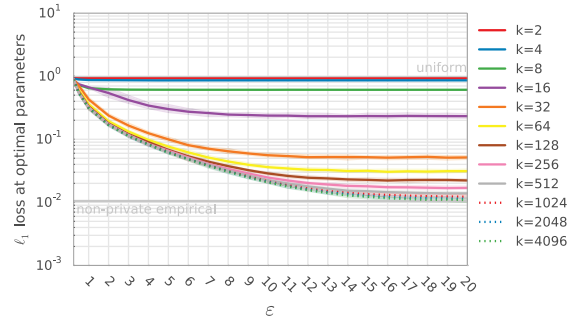


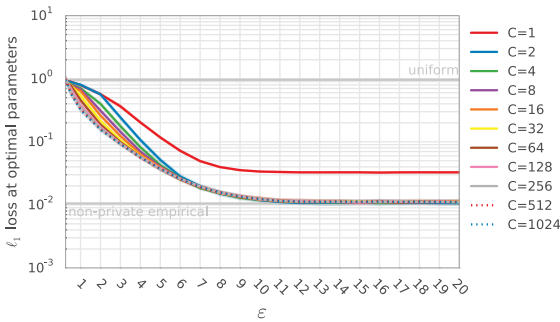
Figure 10: Taking ℓ_1 loss (the utility) and n (the number of users) as fixed requirements (as is the case in many practical scenarios), we approximate the degree of privacy ϵ that can be obtained under O-RR and O-RAPPOR for closed alphabets (lower ϵ is better). Input is generated from an alphabet of $S = 256$ symbols under a geometric distribution with mean= $S/5$, as depicted in Figure 4. Free parameters are set via grid search to minimize the median loss over 50 samples at the given ϵ and fixed parameter values.



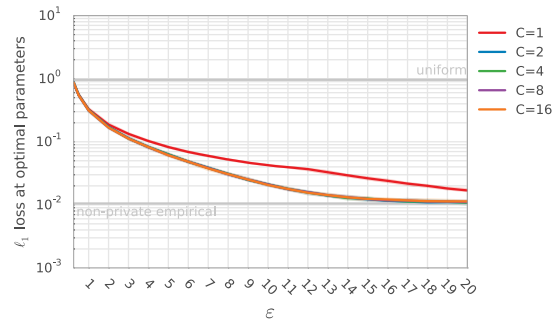
(a) O-RR varying k



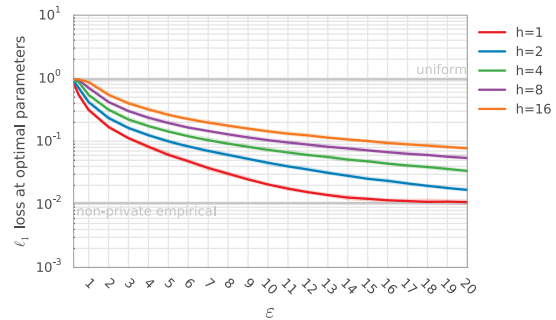
(b) O-RAPPOR varying k



(c) O-RR varying C

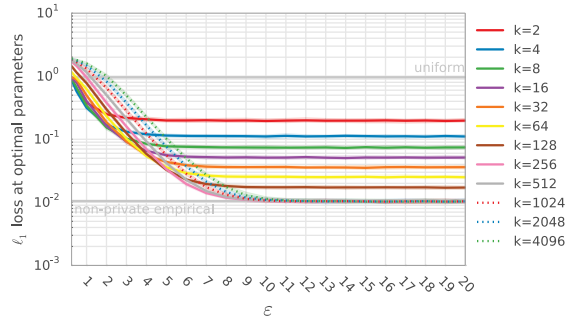


(d) O-RAPPOR varying C

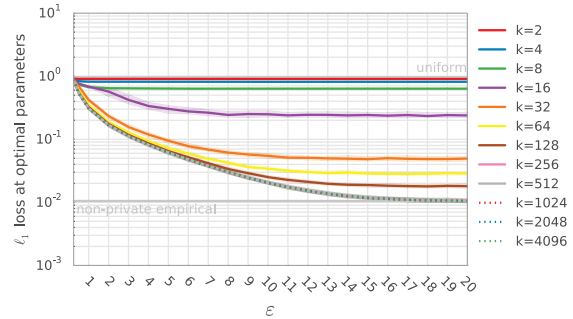


(e) O-RAPPOR varying h

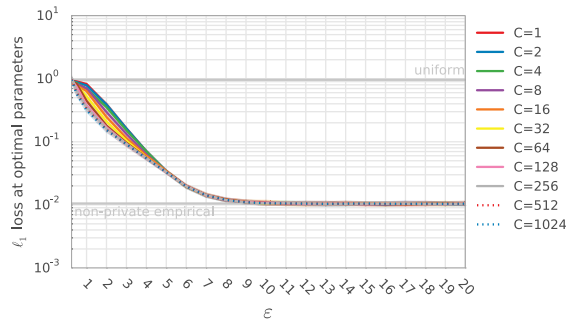
Figure 11: ℓ_1 loss when decoding open alphabets using O-RR and O-RAPPOR under various parameter settings, for $n = 10^6$ users with input drawn from an alphabet of $S = 4096$ symbols under a geometric distribution with mean= $S/5$. Remaining free parameters are set via grid search to minimize the median loss over 50 samples at the given ϵ and fixed parameter values. Lines show median ℓ_1 loss while the (narrow) shaded regions indicate 90% confidence intervals (over 50 samples for the optimal parameter settings.)



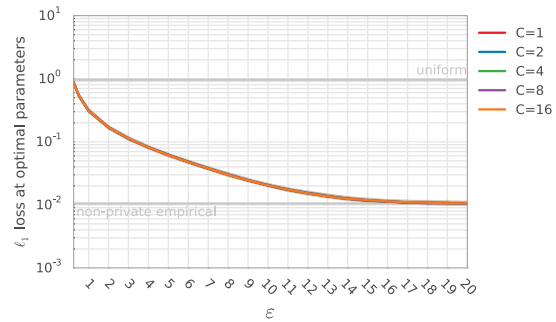
(a) O-RR varying k



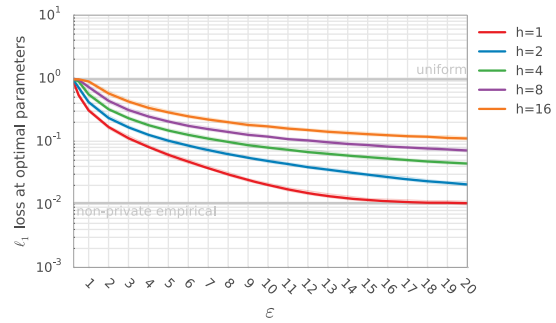
(b) O-RAPPOR varying k



(c) O-RR varying C



(d) O-RAPPOR varying C



(e) O-RAPPOR varying h

Figure 12: ℓ_1 loss when decoding closed alphabets using the O-RR and O-RAPPOR under various parameter settings, for $n = 10^6$ users with input drawn from an alphabet of $S = 4096$ symbols under a geometric distribution with mean= $S/5$. Remaining free parameters are set via grid search to minimize the median loss over 50 samples at the given ϵ and fixed parameter values. Lines show median ℓ_1 loss while the (narrow) shaded regions indicate 90% confidence intervals (over 50 samples for the optimal parameter settings.)

Discrete Distribution Estimation under Local Privacy

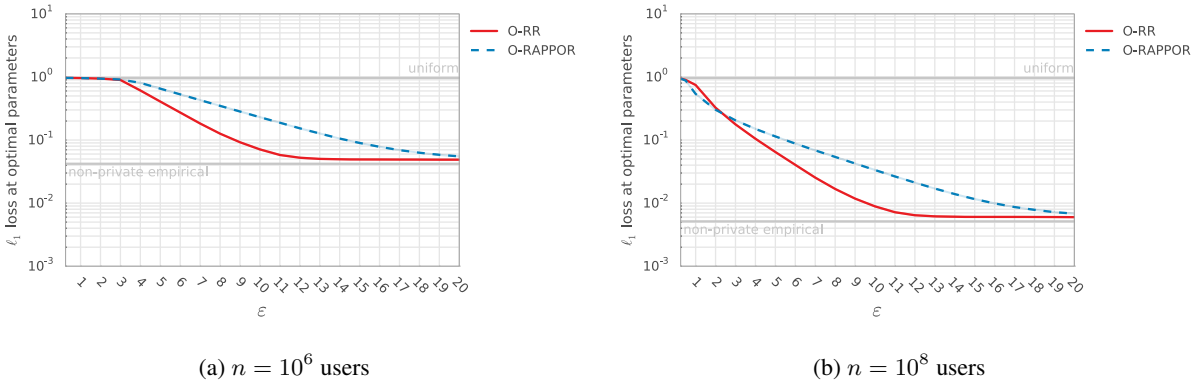


Figure 13: ℓ_1 loss when decoding open alphabets using the O-RR and O-RAPPOR, with input drawn from an alphabet of $S = 4096$ symbols under a geometric distribution with mean= $S/5$. Free parameters are set via grid search over $k \in [2, 4, 8, \dots, 8192, 16384]$, $c \in [1, 2, 4, \dots, 512, 1024]$, $h \in [1, 2]$ to minimize the median loss over 50 samples at the given ϵ value. Lines show median ℓ_1 loss while the (narrow) shaded regions indicate 90% confidence intervals (over 50 samples). Baselines indicate expected loss from (1) using an empirical estimator directly on the input s and (2) using the uniform distribution as the \hat{p} estimate.

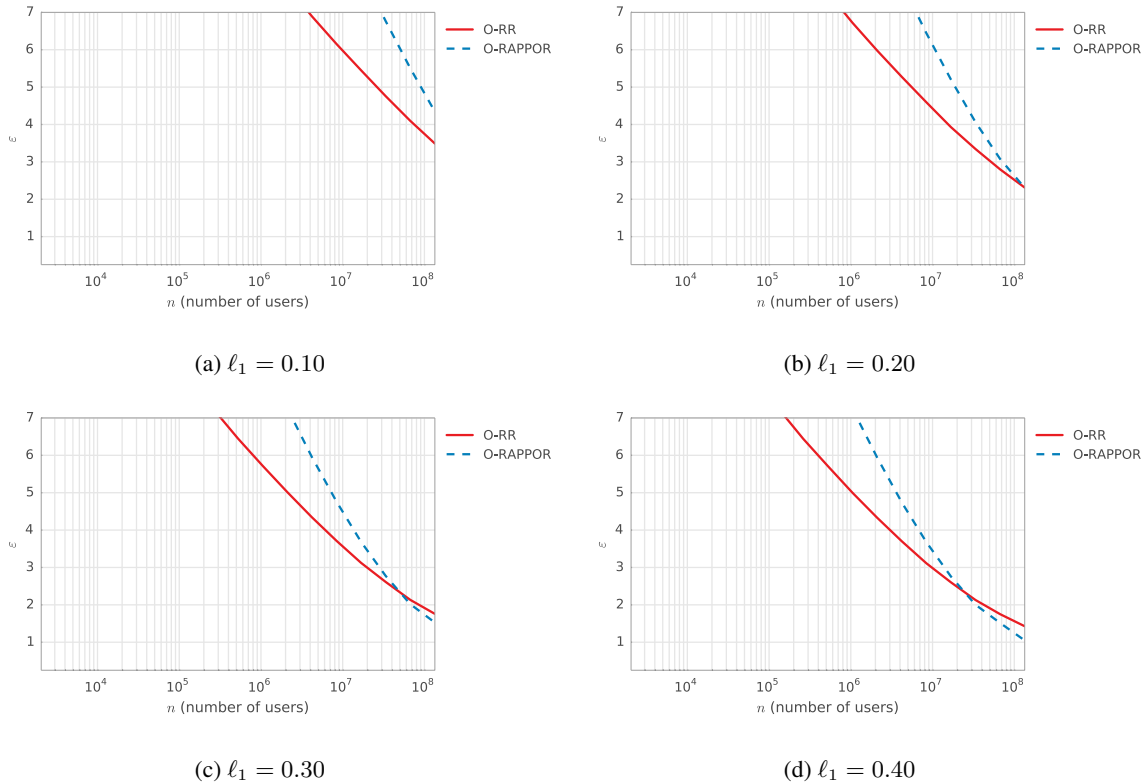
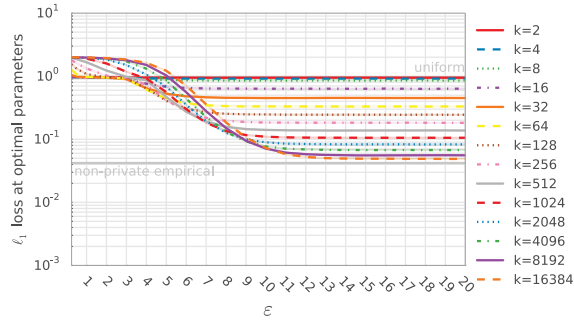
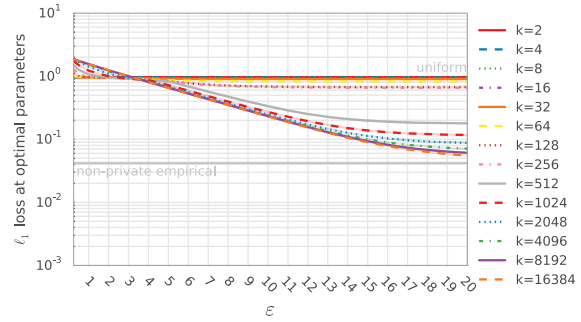


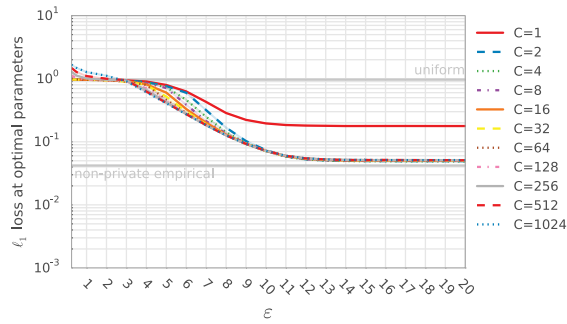
Figure 14: Taking ℓ_1 loss (the utility) and n (the number of users) as fixed requirements (as is the case in many practical scenarios), we approximate the degree of privacy ϵ that can be obtained under O-RR and O-RAPPOR for open alphabets (lower ϵ is better). Input is generated from an alphabet of $S = 4096$ symbols under a geometric distribution with mean= $S/5$, as depicted in Figure 4. Free parameters are set via grid search to minimize the median loss over 50 samples at the given ϵ and fixed parameter values.



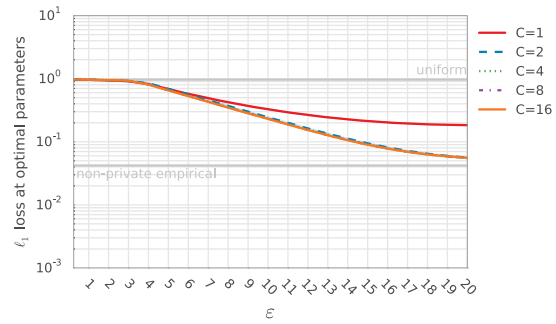
(a) O-RR varying k



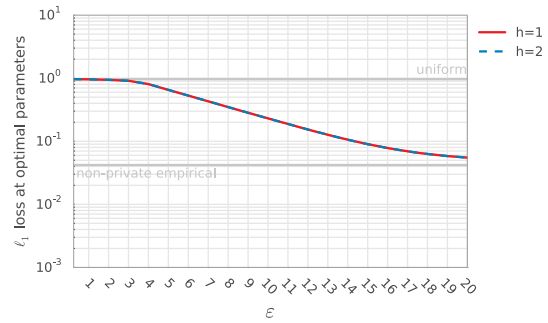
(b) O-RAPPOR varying k



(c) O-RR varying C

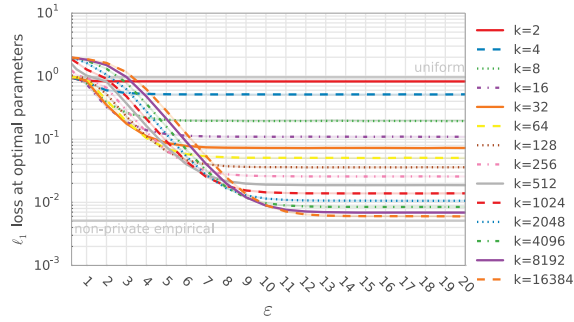


(d) O-RAPPOR varying C

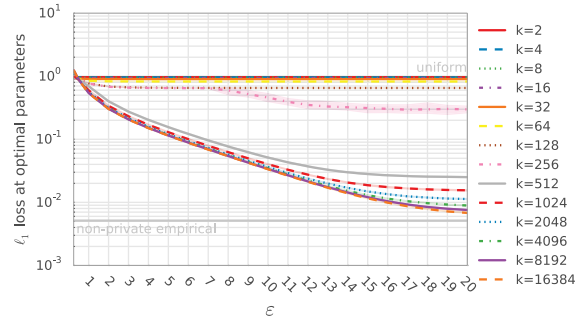


(e) O-RAPPOR varying h

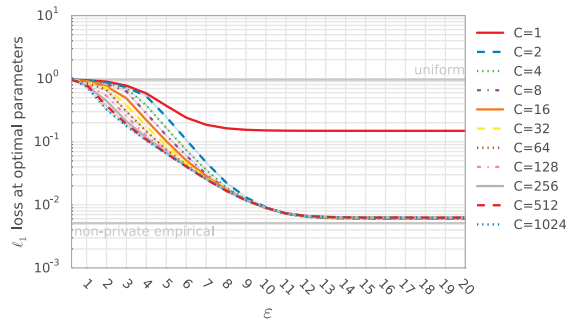
Figure 15: ℓ_1 loss when decoding open alphabets using O-RR and O-RAPPOR under various parameter settings, for $n = 10^6$ users with input drawn from an alphabet of $S = 4096$ symbols under a geometric distribution with mean= $S/5$. Remaining free parameters are set via grid search to minimize the median loss over 50 samples at the given ϵ and fixed parameter values. Lines show median ℓ_1 loss while the (narrow) shaded regions indicate 90% confidence intervals (over 50 samples for the optimal parameter settings.)



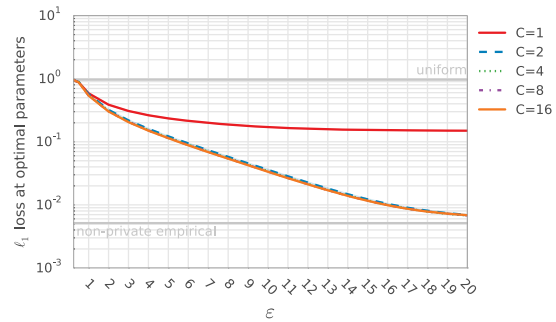
(a) O-RR varying k



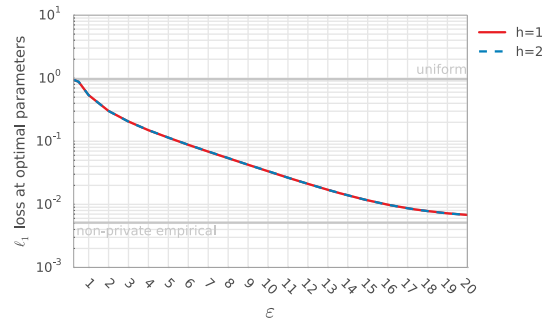
(b) O-RAPPOR varying k



(c) O-RR varying C



(d) O-RAPPOR varying C



(e) O-RAPPOR varying h

Figure 16: ℓ_1 loss when decoding open alphabets using O-RR and O-RAPPOR under various parameter settings, for $n = 10^8$ users with input drawn from an alphabet of $S = 4096$ symbols under a geometric distribution with mean= $S/5$. Remaining free parameters are set via grid search to minimize the median loss over 50 samples at the given ϵ and fixed parameter values. Lines show median ℓ_1 loss while the (narrow) shaded regions indicate 90% confidence intervals (over 50 samples for the optimal parameter settings.)