

ACCURACY FIRST: SELECTING A DIFFERENTIAL PRIVACY LEVEL FOR ACCURACY-CONSTRAINED ERM

KATRINA LIGETT, SETH NEEL, AARON ROTH, BO WAGGONER, AND ZHIWEI STEVEN WU

Hebrew University

University of Pennsylvania

University of Pennsylvania

Microsoft Research

University of Minnesota

e-mail address: zsw@umn.edu

ABSTRACT. Traditional approaches to differential privacy assume a fixed privacy requirement ϵ for a computation, and attempt to maximize the accuracy of the computation subject to the privacy constraint. As differential privacy is increasingly deployed in practical settings, it may often be that there is instead a fixed accuracy requirement for a given computation and the data analyst would like to maximize the privacy of the computation subject to the accuracy constraint. This raises the question of how to find and run a maximally private empirical risk minimizer subject to a given accuracy requirement. We propose a general “noise reduction” framework that can apply to a variety of private empirical risk minimization (ERM) algorithms, using them to “search” the space of privacy levels to find the empirically strongest one that meets the accuracy constraint, and incurring only logarithmic overhead in the number of privacy levels searched. The privacy analysis of our algorithm leads naturally to a version of differential privacy where the privacy parameters are dependent on the data, which we term *ex-post* privacy, and which is related to the recently introduced notion of privacy odometers. We also give an *ex-post* privacy analysis of the classical AboveThreshold privacy tool, modifying it to allow for queries chosen depending on the database. Finally, we apply our approach to two common objective functions, regularized linear and logistic regression, and empirically compare our noise reduction methods to (i) inverting the theoretical utility guarantees of standard private ERM algorithms and (ii) a stronger, empirical baseline based on binary search.

Key words and phrases: differential privacy, empirical risk minimization, accuracy first.

1. INTRODUCTION AND RELATED WORK

Differential Privacy (7; 8) enjoys over a decade of study as a theoretical construct, and a much more recent set of large-scale practical deployments, including by Google (10) and Apple (11). As the large theoretical literature is put into practice, we start to see disconnects between assumptions implicit in the theory and the practical necessities of applications. In this paper we focus our attention on one such assumption in the domain of private empirical risk minimization (ERM): that the data analyst first chooses a privacy requirement, and then attempts to obtain the best accuracy guarantee (or empirical performance) that she can, given the chosen privacy constraint. Existing theory is tailored to this view: the data analyst can pick her privacy parameter ε via some exogenous process, and either plug it into a “utility theorem” to upper bound her accuracy loss, or simply deploy her algorithm and (privately) evaluate its performance. There is a rich and substantial literature on private convex ERM that takes this approach, weaving tight connections between standard mechanisms in differential privacy and standard tools for empirical risk minimization. These methods for private ERM include output and objective perturbation (5; 14; 18; 4), covariance perturbation (19), the exponential mechanism (16; 2), and stochastic gradient descent (2; 21; 12; 6; 20).

While these existing algorithms take a privacy-first perspective, in practice, product requirements may impose hard accuracy constraints, and privacy (while desirable) may not be the over-riding concern. In such situations, things are reversed: the data analyst first fixes an accuracy requirement, and then would like to find the smallest privacy parameter consistent with the accuracy constraint. Here, we find a gap between theory and practice. The only theoretically sound method available is to take a “utility theorem” for an existing private ERM algorithm and solve for the smallest value of ε (the differential privacy parameter)—and other parameter values that need to be set—consistent with her accuracy requirement, and then run the private ERM algorithm with the resulting ε . But because utility theorems tend to be worst-case bounds, this approach will generally be extremely conservative, leading to a much larger value of ε (and hence a much larger leakage of information) than is necessary for the problem at hand. Alternately, the analyst could attempt an empirical search for the smallest value of ε consistent with her accuracy goals. However, because this search is itself a data-dependent computation, it incurs the overhead of additional privacy loss. Furthermore, it is not *a priori* clear how to undertake such a search with nontrivial privacy guarantees for two reasons: first, the worst case could involve a very long search which reveals a large amount of information, and second, the selected privacy parameter is now itself a data-dependent quantity, and so it is not sensible to claim a “standard” guarantee of differential privacy for any finite value of ε ex-ante.

In this paper, we provide a principled variant of this second approach, which attempts to empirically find the smallest value of ε consistent with an accuracy requirement. We give a meta-method that can be applied to several interesting classes of private learning algorithms and introduces very little privacy overhead as a result of the privacy-parameter search. Conceptually, our meta-method initially computes a very private hypothesis, and then gradually subtracts noise (making the computation less and less private) until a sufficient level of accuracy is achieved. One key technique that significantly reduces privacy loss over naive search is the use of correlated noise generated by the method of (15), which formalizes the conceptual idea of “subtracting” noise without incurring additional privacy overhead. In order to select the most private of these queries that meets the accuracy requirement, we introduce a natural modification of the now-classic AboveThreshold algorithm (8), which iteratively

checks a sequence of queries on a dataset and privately releases the index of the first to approximately exceed some fixed threshold. Its privacy cost increases only logarithmically with the number of queries. We provide an analysis of AboveThreshold that holds even if the queries themselves are the result of differentially private computations, showing that if AboveThreshold terminates after t queries, one only pays the privacy costs of AboveThreshold plus the privacy cost of revealing those first t private queries. When combined with the above-mentioned correlated noise technique of (15), this gives an algorithm whose privacy loss is *equal* to that of the final hypothesis output – the previous ones coming “for free” – plus the privacy loss of AboveThreshold. Because the privacy guarantees achieved by this approach are not fixed *a priori*, but rather are a function of the data, we introduce and apply a new, corresponding privacy notion, which we term *ex-post* privacy, and which is closely related to the recently introduced notion of “privacy odometers” (17).

In Section 4, we empirically evaluate our noise reduction meta-method, which applies to any ERM technique which can be described as a post-processing of the Laplace mechanism. This includes both direct applications of the Laplace mechanism, like *output perturbation* (5); and more sophisticated methods like *covariance perturbation* (19), which perturbs the covariance matrix of the data and then performs an optimization using the noisy data. Our experiments concentrate on ℓ_2 regularized least-squares regression and ℓ_2 regularized logistic regression, and we apply our noise reduction meta-method to both output perturbation and covariance perturbation. Our empirical results show that the active, ex-post privacy approach massively outperforms inverting the theory curve, and also improves on a baseline “ ε -doubling” approach.

2. PRIVACY BACKGROUND AND TOOLS

2.1. Differential Privacy and Ex-Post Privacy. Let \mathcal{X} denote the data domain. We call two *datasets* $D, D' \in \mathcal{X}^*$ *neighbors* (written as $D \sim D'$) if D can be derived from D' by replacing a single data point with some other element of \mathcal{X} .

Definition 2.1 (Differential Privacy (7)). Fix $\varepsilon \geq 0$. A randomized algorithm $A : \mathcal{X}^* \rightarrow \mathcal{O}$ is ε -differentially private if for every pair of neighboring data sets $D \sim D' \in \mathcal{X}^*$, and for every event $S \subseteq \mathcal{O}$:

$$\Pr[A(D) \in S] \leq \exp(\varepsilon) \Pr[A(D') \in S].$$

We call $\exp(\varepsilon)$ the *privacy risk* factor.

It is possible to design computations that do not satisfy the differential privacy definition, but whose outputs are private to an extent that can be quantified after the computation halts. For example, consider an experiment that repeatedly runs an ε' -differentially private algorithm, until a stopping condition defined by the output of the algorithm itself is met. This experiment does not satisfy ε -differential privacy for any fixed value of ε , since there is no fixed maximum number of rounds for which the experiment will run (for a fixed number of rounds, a simple composition theorem, Theorem 2.5, shows that the ε -guarantees in a sequence of computations “add up.”) However, if ex-post we see that the experiment has stopped after k rounds, the data can in some sense be assured an “ex-post privacy loss” of only $k\varepsilon'$. Rogers et al. (17) initiated the study of *privacy odometers*, which formalize this idea. They study privacy composition when the data analyst can choose the privacy parameters of subsequent computations as a function of the outcomes of previous computations.

We apply a related idea here, for a different purpose. Our goal is to design one-shot algorithms that always achieve a target accuracy but that may have variable privacy levels depending on their input.

Definition 2.2. Given a randomized algorithm $\mathcal{A} : \mathcal{X}^* \rightarrow \mathcal{O}$, define the *ex-post privacy loss*¹ of \mathcal{A} on outcome o to be

$$\text{Loss}(o) = \max_{D, D' : D \sim D'} \log \frac{\Pr[\mathcal{A}(D) = o]}{\Pr[\mathcal{A}(D') = o]}.$$

We refer to $\exp(\text{Loss}(o))$ as the *ex-post privacy risk factor*.

Definition 2.3 (Ex-Post Differential Privacy). Let $\mathcal{E} : \mathcal{O} \rightarrow (\mathbb{R}_{\geq 0} \cup \{\infty\})$ be a function on the outcome space of algorithm $\mathcal{A} : \mathcal{X}^* \rightarrow \mathcal{O}$. Given an outcome $o = \mathcal{A}(D)$, we say that \mathcal{A} satisfies $\mathcal{E}(o)$ -*ex-post* differential privacy if for all $o \in \mathcal{O}$, $\text{Loss}(o) \leq \mathcal{E}(o)$.

Note that if $\mathcal{E}(o) \leq \varepsilon$ for all o , \mathcal{A} is ε -differentially private. Ex-post differential privacy has the same semantics as differential privacy, once the output of the mechanism is known: it bounds the log-likelihood ratio of the dataset being D vs. D' , which controls how an adversary with an arbitrary prior on the two cases can update her posterior.

2.2. Differential Privacy Tools. Differentially private computations enjoy two nice properties:

Theorem 2.4 Post Processing (7). *Let $A : \mathcal{X}^* \rightarrow \mathcal{O}$ be any ε -differentially private algorithm, and let $f : \mathcal{O} \rightarrow \mathcal{O}'$ be any function. Then the algorithm $f \circ A : \mathcal{X}^* \rightarrow \mathcal{O}'$ is also ε -differentially private.*

Post-processing implies that, for example, every *decision* process based on the output of a differentially private algorithm is also differentially private.

Theorem 2.5 Composition (7). *Let $A_1 : \mathcal{X}^* \rightarrow \mathcal{O}$, $A_2 : \mathcal{X}^* \rightarrow \mathcal{O}'$ be algorithms that are ε_1 - and ε_2 -differentially private, respectively. Then the algorithm $A : \mathcal{X}^* \rightarrow \mathcal{O} \times \mathcal{O}'$ defined as $A(x) = (A_1(x), A_2(x))$ is $(\varepsilon_1 + \varepsilon_2)$ -differentially private.*

The composition theorem holds even if the composition is *adaptive*—see (9) for details.

The Laplace mechanism. The most basic subroutine we will use is the *Laplace mechanism*. The Laplace Distribution centered at 0 with scale b is the distribution with probability density function $\text{Lap}(z|b) = \frac{1}{2b} e^{-\frac{|z|}{b}}$. We say $X \sim \text{Lap}(b)$ when X has Laplace distribution with scale b . Let $f : \mathcal{X}^* \rightarrow \mathbb{R}^d$ be an arbitrary d -dimensional function. The ℓ_1 sensitivity of f is defined to be $\Delta_1(f) = \max_{D \sim D'} \|f(D) - f(D')\|_1$. The *Laplace mechanism* with parameter ε simply adds noise drawn independently from $\text{Lap}\left(\frac{\Delta_1(f)}{\varepsilon}\right)$ to each coordinate of $f(x)$.

Theorem 2.6 (7). *The Laplace mechanism is ε -differentially private.*

Gradual private release. Koufogiannis et al. (15) study how to gradually release private data using the Laplace mechanism with an increasing sequence of ε values, with a privacy cost scaling only with the privacy of the *marginal* distribution on the least private release, rather than the sum of the privacy costs of independent releases. For intuition, the

¹If \mathcal{A} 's output is from a continuous distribution rather than discrete, we abuse notation and write $\Pr[\mathcal{A}(D) = o]$ to mean the probability density at output o .

algorithm can be pictured as a continuous random walk starting at some private data v with the property that the marginal distribution at each point in time is Laplace centered at v , with variance increasing over time. Releasing the value of the random walk at a fixed point in time gives a certain output distribution, for example, \hat{v} , with a certain privacy guarantee ε . To produce \hat{v}' whose *ex-ante* distribution has higher variance (is more private), one can simply “fast forward” the random walk from a starting point of \hat{v} to reach \hat{v}' ; to produce a less private \hat{v}' , one can “rewind.” The total privacy cost is $\max\{\varepsilon, \varepsilon'\}$ because, given the “least private” point (say \hat{v}), all “more private” points can be derived as post-processings given by taking a random walk of a certain length starting at \hat{v} . Note that were the Laplace random variables used for each release independent, the composition theorem would require *summing* the ε values of all releases.

In our private algorithms, we will use their noise reduction mechanism as a building block to generate a list of private hypotheses $\theta^1, \dots, \theta^T$ with gradually increasing ε values. Importantly, releasing any prefix $(\theta^1, \dots, \theta^t)$ only incurs the privacy loss in θ^t . More formally:

Algorithm 1 Noise Reduction (15): $\text{NR}(v, \Delta, \{\varepsilon_t\})$

Input: private vector v , sensitivity parameter Δ , list $\varepsilon_1 < \varepsilon_2 < \dots < \varepsilon_T$
 Set $\hat{v}_T := v + \text{Lap}(\Delta/\varepsilon_T)$
for $t = T - 1, T - 2, \dots, 1$ **do**
 With probability $\left(\frac{\varepsilon_t}{\varepsilon_{t+1}}\right)^2$: set $\hat{v}_t := \hat{v}_{t+1}$
 Else: set $\hat{v}_t := \hat{v}_{t+1} + \text{Lap}(\Delta/\varepsilon_t)$
 Return $\hat{v}_1, \dots, \hat{v}_T$

Theorem 2.7 (15). *Let f have ℓ_1 sensitivity Δ and let $\hat{v}_1, \dots, \hat{v}_T$ be the output of Algorithm 1 on $v = f(D)$, Δ , and the increasing list $\varepsilon_1, \dots, \varepsilon_T$. Then for any t , the algorithm which outputs the prefix $(\hat{v}_1, \dots, \hat{v}_t)$ is ε_t -differentially private.*

2.3. AboveThreshold with Private Queries. Our high-level approach to our eventual ERM problem will be as follows: Generate a sequence of hypotheses $\theta_1, \dots, \theta_T$, each with increasing accuracy and decreasing privacy; then test their accuracy levels sequentially, outputting the first one whose accuracy is “good enough.” The classical AboveThreshold algorithm (8) takes in a dataset and a sequence of queries and privately outputs the index of the first query to exceed a given threshold (with some error due to noise). We would like to use AboveThreshold to perform these accuracy checks, but there is an important obstacle: for us, the “queries” themselves depend on the private data.² A standard composition analysis would involve first privately publishing *all* the queries, then running AboveThreshold on these queries (which are now public). Intuitively, though, it would be much better to generate and publish the queries one at a time, until AboveThreshold halts, at which point one would not publish any more queries. The problem with analyzing this approach is that, *a-priori*,

²In fact, there are many applications beyond our own in which the sequence of queries input to AboveThreshold might be the result of some private prior computation on the data, and where we would like to release both the stopping index of AboveThreshold and the “query object.” (In our case, the query objects will be parameterized by learned hypotheses $\theta_1, \dots, \theta_T$.)

we do not know when AboveThreshold will terminate; to address this, we analyze the *ex-post privacy* guarantee of the algorithm.³

Algorithm 2 InteractiveAboveThreshold: $\text{IAT}(D, \varepsilon, W, \Delta, M)$

Input: Dataset D , privacy loss ε , threshold W , ℓ_1 sensitivity Δ , algorithm M
 Let $\hat{W} = W + \text{Lap}\left(\frac{2\Delta}{\varepsilon}\right)$
for each query $t = 1, \dots, T$ **do**
 Query $f_t \leftarrow M(D)_t$
 if $f_t(D) + \text{Lap}\left(\frac{4\Delta}{\varepsilon}\right) \geq \hat{W}$ **then** Output (t, f_t) ; **Halt**.
 Output (T, \perp) .

Let us say that an algorithm $M(D) = (f_1, \dots, f_T)$ is $(\varepsilon_1, \dots, \varepsilon_T)$ -*prefix-private* if for each t , the function that runs $M(D)$ and outputs just the prefix (f_1, \dots, f_t) is ε_t -differentially private.

Lemma 2.8. *Let $M : \mathcal{X}^* \rightarrow (\mathcal{X}^* \rightarrow \mathcal{O})^T$ be a $(\varepsilon_1, \dots, \varepsilon_T)$ -prefix private algorithm that returns T queries, and let each query output by M have ℓ_1 sensitivity at most Δ . Then Algorithm 2 run on $D, \varepsilon_A, W, \Delta$, and M is \mathcal{E} -ex-post differentially private for $\mathcal{E}((t, \cdot)) = \varepsilon_A + \varepsilon_t$ for any $t \in [T]$.*

The proof, which is a variant on the proof of privacy for AboveThreshold (8), appears in the appendix, along with an accuracy theorem for IAT.

3. NOISE-REDUCTION WITH PRIVATE ERM

In this section, we provide a general private ERM framework that allows us to approach the best privacy guarantee achievable on the data given a target excess risk goal. Throughout the section, we consider an input dataset D that consists of n row vectors $X_1, X_2, \dots, X_n \in \mathbb{R}^p$ and a column $y \in \mathbb{R}^n$. We will assume that each $\|X_i\|_1 \leq 1$ and $|y_i| \leq 1$. Let $d_i = (X_i, y_i) \in \mathbb{R}^{p+1}$ be the i -th data record. Let ℓ be a loss function such that for any hypothesis θ and any data point (X_i, y_i) the loss is $\ell(\theta, (X_i, y_i))$. Given an input dataset D and a regularization parameter λ , the goal is to minimize the following regularized empirical loss function over some feasible set C :

$$L(\theta, D) = \frac{1}{n} \sum_{i=1}^n \ell(\theta, (X_i, y_i)) + \frac{\lambda}{2} \|\theta\|_2^2.$$

Let $\theta^* = \arg\min_{\theta \in C} \ell(\theta, D)$. Given a target accuracy parameter α , we wish to privately compute a θ_p that satisfies $L(\theta_p, D) \leq L(\theta^*, D) + \alpha$, while achieving the best ex-post privacy guarantee. For simplicity, we will sometimes write $L(\theta)$ for $L(\theta, D)$.

One simple baseline approach is a “doubling method”: Start with a small ε value, run an ε -differentially private algorithm to compute a hypothesis θ and use the Laplace mechanism to estimate the excess risk of θ ; if the excess risk is lower than the target, output θ ; otherwise double the value of ε and repeat the same process. (See the appendix for details.) As a result, we pay for privacy loss for every hypothesis we compute and every excess risk we estimate.

In comparison, our meta-method provides a more cost-effective way to select the privacy level. The algorithm takes a more refined set of privacy levels $\varepsilon_1 < \dots < \varepsilon_T$ as input and

³This result does not follow from a straightforward application of privacy odometers from (17), because the privacy analysis of algorithms like the noise reduction technique is not compositional.

generates a sequence of hypotheses $\theta^1, \dots, \theta^T$ such that the generation of each θ^t is ε_t -private. Then it releases the hypotheses θ^t in order, halting as soon as a released hypothesis meets the accuracy goal. Importantly, there are two key components that reduce the privacy loss in our method:

- (1) We use Algorithm 1, the “noise reduction” method of (15), for generating the sequence of hypotheses: we first compute a very private and noisy θ^1 , and then obtain the subsequent hypotheses by gradually “de-noising” θ^1 . As a result, any prefix $(\theta^1, \dots, \theta^k)$ incurs a privacy loss of only ε_k (as opposed to $(\varepsilon_1 + \dots + \varepsilon_k)$ if the hypotheses were independent).
- (2) When evaluating the excess risk of each hypothesis, we use Algorithm 2, Interactive-AboveThreshold, to determine if its excess risk exceeds the target threshold. This incurs substantially less privacy loss than independently evaluating the excess risk of each hypothesis using the Laplace mechanism (and hence allows us to search a finer grid of values).

For the rest of this section, we will instantiate our method concretely for two ERM problems: ridge regression and logistic regression. In particular, our noise-reduction method is based on two private ERM algorithms: the recently introduced covariance perturbation technique (19) and the output perturbation method (5).

3.1. Covariance Perturbation for Ridge Regression. In ridge regression, we consider the squared loss function: $\ell((X_i, y_i), \theta) = \frac{1}{2}(y_i - \langle \theta, X_i \rangle)^2$, and hence empirical loss over the data set is defined as

$$L(\theta, D) = \frac{1}{2n} \|y - X\theta\|_2^2 + \frac{\lambda \|\theta\|_2^2}{2},$$

where X denotes the $(n \times p)$ matrix with row vectors X_1, \dots, X_n and $y = (y_1, \dots, y_n)$. Since the optimal solution for the unconstrained problem has ℓ_2 norm no more than $\sqrt{1/\lambda}$ (see the appendix for a proof), we will focus on optimizing θ over the constrained set $C = \{a \in \mathbb{R}^p \mid \|a\|_2 \leq \sqrt{1/\lambda}\}$, which will be useful for bounding the ℓ_1 sensitivity of the empirical loss.

Before we formally introduce the covariance perturbation algorithm due to (19), observe that the optimal solution θ^* can be computed as

$$\theta^* = \operatorname{argmin}_{\theta \in C} L(\theta, D) = \operatorname{argmin}_{\theta \in C} \frac{(\theta^\top (X^\top X) \theta - 2\langle X^\top y, \theta \rangle)}{2n} + \frac{\lambda \|\theta\|_2^2}{2}.$$

In other words, θ^* only depends on the private data through $X^\top y$ and $X^\top X$. To compute a private hypothesis, the covariance perturbation method simply adds Laplace noise to each entry of $X^\top y$ and $X^\top X$ (the covariance matrix), and solves the optimization based on the noisy matrix and vector. The formal description of the algorithm and its guarantee are in Theorem 3.1. Our analysis differs from the one in (19) in that their paper considers the “local privacy” setting, and also adds Gaussian noise whereas we use Laplace. The proof is deferred to the appendix.

Theorem 3.1 . *Fix any $\varepsilon > 0$. For any input data set D , consider the mechanism \mathcal{M} that computes*

$$\theta_p = \operatorname{argmin}_{\theta \in C} \frac{1}{2n} (\theta^\top (X^\top X + B) \theta - 2\langle X^\top y + b, \theta \rangle) + \frac{\lambda \|\theta\|_2^2}{2},$$

where $B \in \mathbb{R}^{p \times p}$ and $b \in \mathbb{R}^{p \times 1}$ are random Laplace matrices such that each entry of B and b is drawn from $\text{Lap}(4/\varepsilon)$. Then \mathcal{M} satisfies ε -differential privacy and the output θ_p satisfies

$$\mathbb{E}_{B,b} [L(\theta_p) - L(\theta^*)] \leq \frac{4\sqrt{2}(2\sqrt{p/\lambda} + p/\lambda)}{n\varepsilon}.$$

In our algorithm COVNR, we will apply the noise reduction method, Algorithm 1, to produce a sequence of noisy versions of the private data $(X^\top X, X^\top y)$: $(Z^1, z^1), \dots, (Z^T, z^T)$, one for each privacy level. Then for each (Z^t, z^t) , we will compute the private hypothesis by solving the noisy version of the optimization problem in Equation (3.1). The full description of our algorithm COVNR is in Algorithm 3, and satisfies the following guarantee:

Algorithm 3 Covariance Perturbation with Noise-Reduction: $\text{COVNR}(D, \{\varepsilon_1, \dots, \varepsilon_T\}, \alpha, \gamma)$

Input: private data set $D = (X, y)$, accuracy parameter α , privacy levels $\varepsilon_1 < \varepsilon_2 < \dots < \varepsilon_T$, and failure probability γ

Instantiate `InteractiveAboveThreshold`: $\mathcal{A} = \text{IAT}(D, \varepsilon_0, -\alpha/2, \Delta, \cdot)$ with $\varepsilon_0 = 16\Delta(\log(2T/\gamma))/\alpha$ and $\Delta = (\sqrt{1/\lambda} + 1)^2/(n)$

Let $C = \{a \in \mathbb{R}^p \mid \|a\|_2 \leq \sqrt{1/\lambda}\}$ and $\theta^* = \text{argmin}_{\theta \in C} L(\theta)$

Compute noisy data:

$$\{Z^t\} = \text{NR}((X^\top X), 2, \{\varepsilon_1/2, \dots, \varepsilon_T/2\}), \quad \{z^t\} = \text{NR}((X^\top Y), 2, \{\varepsilon_1/2, \dots, \varepsilon_T/2\})$$

for $t = 1, \dots, T$: **do**

$$\theta^t = \text{argmin}_{\theta \in C} \frac{1}{2n} (\theta^\top Z^t \theta - 2\langle z^t, \theta \rangle) + \frac{\lambda \|\theta\|_2^2}{2} \quad (3.1)$$

Let $f^t(D) = L(\theta^*, D) - L(\theta^t, D)$; Query \mathcal{A} with query f^t to check accuracy

if \mathcal{A} returns (t, f^t) **then Output** (t, θ^t)

Output: (\perp, θ^*)

Theorem 3.2 . *The instantiation of $\text{COVNR}(D, \{\varepsilon_1, \dots, \varepsilon_T\}, \alpha, \gamma)$ outputs a hypothesis θ_p that with probability $1 - \gamma$ satisfies $L(\theta_p) - L(\theta^*) \leq \alpha$. Moreover, it is \mathcal{E} -ex-post differentially private, where the privacy loss function $\mathcal{E}: ([T] \cup \{\perp\}) \times \mathbb{R}^p \rightarrow (\mathbb{R}_{\geq 0} \cup \{\infty\})$ is defined as $\mathcal{E}((k, \cdot)) = \varepsilon_0 + \varepsilon_k$ for any $k \neq \perp$, $\mathcal{E}((\perp, \cdot)) = \infty$, and*

$$\varepsilon_0 = \frac{16(\sqrt{1/\lambda} + 1)^2 \log(2T/\gamma)}{n\alpha}$$

is the privacy loss incurred by IAT.

3.2. Output Perturbation for Logistic Regression. Next, we show how to combine the output perturbation method with noise reduction for the ridge regression problem.⁴ In this setting, the input data consists of n labeled examples $(X_1, y_1), \dots, (X_n, y_n)$, such that for each i , $X_i \in \mathbb{R}^p$, $\|X_i\|_1 \leq 1$, and $y_i \in \{-1, 1\}$. The goal is to train a linear classifier given

⁴We study the ridge regression problem for concreteness. Our method works for any ERM problem with strongly convex loss functions.

by a weight vector θ for the examples from the two classes. We consider the logistic loss function: $\ell(\theta, (X_i, y_i)) = \log(1 + \exp(-y_i \theta^\top X_i))$, and the empirical loss is

$$L(\theta, D) = \frac{1}{n} \sum_{i=1}^n \log(1 + \exp(-y_i \theta^\top X_i)) + \frac{\lambda \|\theta\|_2^2}{2}.$$

The output perturbation method simply adds Laplace noise to perturb each coordinate of the optimal solution θ^* . The following is the formal guarantee of output perturbation. Our analysis deviates slightly from the one in (5) since we are adding Laplace noise (see the appendix).

Theorem 3.3 . *Fix any $\varepsilon > 0$. Let $r = \frac{2\sqrt{p}}{n\lambda\varepsilon}$. For any input dataset D , consider the mechanism that first computes $\theta^* = \operatorname{argmin}_{\theta \in \mathbb{R}^p} L(\theta)$, then outputs $\theta_p = \theta^* + b$, where b is a random vector with its entries drawn i.i.d. from $\text{Lap}(r)$. Then \mathcal{M} satisfies ε -differential privacy, and θ_p has excess risk*

$$\mathbb{E}_b [L(\theta_p) - L(\theta^*)] \leq \frac{2\sqrt{2}p}{n\lambda\varepsilon} + \frac{4p^2}{n^2\lambda\varepsilon^2}.$$

Given the output perturbation method, we can simply apply the noise reduction method NR to the optimal hypothesis θ^* to generate a sequence of noisy hypotheses. We will again use InteractiveAboveThreshold to check the excess risk of the hypotheses. The full algorithm OUTPUTNR follows the same structure in Algorithm 3, and we defer the formal description to the appendix.

Theorem 3.4 . *The instantiation of $\text{OUTPUTNR}(D, \varepsilon_0, \{\varepsilon_1, \dots, \varepsilon_T\}, \alpha, \gamma)$ is \mathcal{E} -ex-post differentially private and outputs a hypothesis θ_p that with probability $1 - \gamma$ satisfies $L(\theta_p) - L(\theta^*) \leq \alpha$, where the privacy loss function $\mathcal{E}: ([T] \cup \{\perp\}) \times \mathbb{R}^p \rightarrow (\mathbb{R}_{\geq 0} \cup \{\infty\})$ is defined as $\mathcal{E}((k, \cdot)) = \varepsilon_0 + \varepsilon_k$ for any $k \neq \perp$, $\mathcal{E}((\perp, \cdot)) = \infty$, and*

$$\varepsilon_0 \leq \frac{32 \log(2T/\gamma) \sqrt{2 \log 2/\lambda}}{n\alpha}$$

is the privacy loss incurred by IAT.

Proof sketch of Theorems 3.2 and 3.4. The accuracy guarantees for both algorithms follow from an accuracy guarantee of the IAT algorithm (a variant on the standard AboveThreshold bound) and the fact that we output θ^* if IAT identifies no accurate hypothesis. For the privacy guarantee, first note that any prefix of the noisy hypotheses $\theta^1, \dots, \theta^t$ satisfies ε_t -differential privacy because of our instantiation of the Laplace mechanism (see the appendix for the ℓ_1 sensitivity analysis) and noise-reduction method NR. Then the ex-post privacy guarantee directly follows Lemma 2.8. \square

4. EXPERIMENTS

To evaluate the methods described above, we conducted empirical evaluations in two settings. We used ridge regression to predict (log) popularity of posts on Twitter in the dataset of (1), with $p = 77$ features and subsampled to $n = 100,000$ data points. Logistic regression was applied to classifying network events as innocent or malicious in the KDD-99 Cup

dataset (13), with 38 features and subsampled to 100,000 points. Details of parameters and methods appear in the appendix.⁵

In each case, we tested the algorithm’s average ex-post privacy loss for a range of input accuracy goals α , fixing a modest failure probability $\gamma = 0.1$ (and we observed that excess risks were concentrated well below $\alpha/2$, suggesting a pessimistic analysis). The results show our meta-method gives a large improvement over the “theory” approach of simply inverting utility theorems for private ERM algorithms. (In fact, the utility theorem for the popular private stochastic gradient descent algorithm does not even give meaningful guarantees for the ranges of parameters tested; one would need an order of magnitude more data points, and even then the privacy losses are enormous, perhaps due to loose constants in the analysis.)

To gauge the more modest improvement over DOUBLINGMETHOD, note that the variation in the privacy risk factor e^ϵ can still be very large; for instance, in the ridge regression setting of $\alpha = 0.05$, Noise Reduction has $e^\epsilon \approx 10.0$ while DOUBLINGMETHOD has $e^\epsilon \approx 495$; at $\alpha = 0.075$, the privacy risk factors are 4.65 and 56.6 respectively.

Interestingly, for our meta-method, the contribution to privacy loss from “testing” hypotheses (the InteractiveAboveThreshold technique) was significantly larger than that from “generating” them (NoiseReduction). One place where the InteractiveAboveThreshold analysis is loose is in using a theoretical bound on the maximum norm of any hypothesis to compute the sensitivity of queries. The actual norms of hypotheses tested was significantly lower which, if taken as guidance to the practitioner in advance, would drastically improve the privacy guarantee of both adaptive methods.

5. FUTURE DIRECTIONS

Throughout this paper, we focus on ϵ -differential privacy, instead of the weaker (ϵ, δ) - (approximate) differential privacy. Part of the reason is that an analogue of Lemma 2.8 does not seem to hold for (ϵ, δ) -differentially private queries without further assumptions, as the necessity to union-bound over the δ “failure probability” that the privacy loss is bounded for each query can erase the ex-post gains. We leave obtaining similar results for approximate differential privacy as an open problem. More generally, we wish to extend our ex-post privacy framework to approximate differential privacy, or to the stronger notion of concentrated differential privacy (3). Such results will allow us to obtain ex-post privacy guarantees for a much broader class of algorithms.

ACKNOWLEDGMENT

The authors wish to acknowledge fruitful discussions with A and B.

REFERENCES

- [1] The AMA Team at Laboratoire d’Informatique de Grenoble. Buzz prediction in on-line social media, 2017. URL: <http://ama.liglab.fr/resourcestools/datasets/buzz-prediction-in-social-media/>.
- [2] Raef Bassily, Adam D. Smith, and Abhradeep Thakurta. Private empirical risk minimization, revisited. *CoRR*, abs/1405.7085, 2014. URL: <http://arxiv.org/abs/1405.7085>.

⁵ A full implementation of our algorithms appears at <https://github.com/steven7woo/Accuracy-First-Differential-Privacy> and (22).

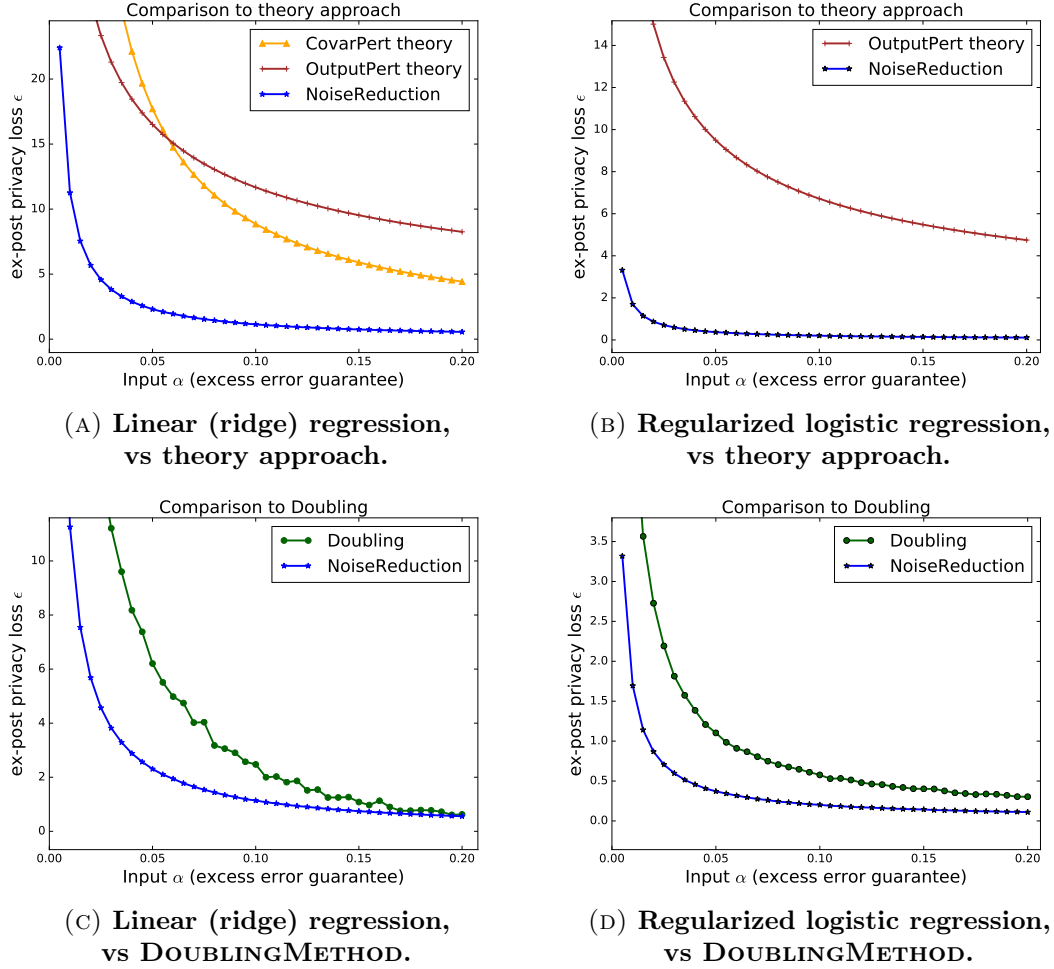


FIGURE 1. **Ex-post privacy loss.** (1a) and (1c), left, represent ridge regression on the Twitter dataset, where Noise Reduction and DOUBLINGMETHOD both use Covariance Perturbation. (1b) and (1d), right, represent logistic regression on the KDD-99 Cup dataset, where both Noise Reduction and DOUBLINGMETHOD use Output Perturbation. The top plots compare Noise Reduction to the “theory approach”: running the algorithm once using the value of ϵ that guarantees the desired expected error via a utility theorem. The bottom compares to the DOUBLINGMETHOD baseline. Note the top plots are generous to the theory approach: the theory curves promise only expected error, whereas Noise Reduction promises a high probability guarantee. Each point is an average of 80 trials (Twitter dataset) or 40 trials (KDD-99 dataset).

- [3] Mark Bun and Thomas Steinke. Concentrated differential privacy: Simplifications, extensions, and lower bounds. In *Theory of Cryptography - 14th International Conference, TCC 2016-B, Beijing, China, October 31 - November 3, 2016, Proceedings, Part I*, pages 635–658, 2016. URL: https://doi.org/10.1007/978-3-662-53641-4_24, doi: 10.1007/978-3-662-53641-4_24.

- [4] Kamalika Chaudhuri and Claire Monteleoni. Privacy-preserving logistic regression. In *Advances in Neural Information Processing Systems 21, Proceedings of the Twenty-Second Annual Conference on Neural Information Processing Systems, Vancouver, British Columbia, Canada, December 8-11, 2008*, pages 289–296, 2008. URL: <http://papers.nips.cc/paper/3486-privacy-preserving-logistic-regression>.
- [5] Kamalika Chaudhuri, Claire Monteleoni, and Anand D. Sarwate. Differentially private empirical risk minimization. *Journal of Machine Learning Research*, 12:1069–1109, 2011. URL: <http://dl.acm.org/citation.cfm?id=2021036>.
- [6] John C. Duchi, Michael I. Jordan, and Martin J. Wainwright. Local privacy and statistical minimax rates. In *51st Annual Allerton Conference on Communication, Control, and Computing, Allerton 2013, Allerton Park & Retreat Center, Monticello, IL, USA, October 2-4, 2013*, page 1592, 2013. URL: <http://dx.doi.org/10.1109/Allerton.2013.6736718>, doi:10.1109/Allerton.2013.6736718.
- [7] Cynthia Dwork, Frank McSherry, Kobbi Nissim, and Adam Smith. Calibrating noise to sensitivity in private data analysis. In *Theory of Cryptography Conference*, pages 265–284. Springer, 2006.
- [8] Cynthia Dwork and Aaron Roth. The algorithmic foundations of differential privacy. *Foundations and Trends in Theoretical Computer Science*, 9(3–4):211–407, 2014.
- [9] Cynthia Dwork, Guy N Rothblum, and Salil Vadhan. Boosting and differential privacy. In *Foundations of Computer Science (FOCS), 2010 51st Annual IEEE Symposium on*, pages 51–60. IEEE, 2010.
- [10] Giulia Fanti, Vasyli Pihur, and Úlfar Erlingsson. Building a RAPPOR with the unknown: Privacy-preserving learning of associations and data dictionaries. *Proceedings on Privacy Enhancing Technologies (PoPETS)*, issue 3, 2016, 2016.
- [11] Andy Greenberg. Apple’s ‘differential privacy’ is about collecting your data—but not your data. *Wired Magazine*, 2016. URL: <https://www.wired.com/2016/06/apples-differential-privacy-collecting-data/>.
- [12] Prateek Jain, Praveesh Kothari, and Abhradeep Thakurta. Differentially private online learning. In *COLT 2012 - The 25th Annual Conference on Learning Theory, June 25-27, 2012, Edinburgh, Scotland*, pages 24.1–24.34, 2012. URL: <http://www.jmlr.org/proceedings/papers/v23/jain12/jain12.pdf>.
- [13] KDD’99. KDD cup 1999 data, 1999. URL: <http://kdd.ics.uci.edu/databases/kddcup99/kddcup99.html>.
- [14] Daniel Kifer, Adam D. Smith, and Abhradeep Thakurta. Private convex optimization for empirical risk minimization with applications to high-dimensional regression. In *COLT 2012 - The 25th Annual Conference on Learning Theory, June 25-27, 2012, Edinburgh, Scotland*, pages 25.1–25.40, 2012. URL: <http://www.jmlr.org/proceedings/papers/v23/kifer12/kifer12.pdf>.
- [15] Fragkiskos Koufogiannis, Shuo Han, and George J. Pappas. Gradual release of sensitive data under differential privacy. *Journal of Privacy and Confidentiality*, 7, 2017.
- [16] 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.
- [17] Ryan M Rogers, Aaron Roth, Jonathan Ullman, and Salil Vadhan. Privacy odometers and filters: Pay-as-you-go composition. In D. D. Lee, M. Sugiyama, U. V. Luxburg, I. Guyon, and R. Garnett, editors, *Advances in Neural Information Processing Systems 29*, pages 1921–1929. Curran Associates, Inc., 2016. URL: <http://papers.nips.cc/>

- [paper/6170-privacy-odometers-and-filters-pay-as-you-go-composition.pdf](#).
- [18] Benjamin I. P. Rubinstein, Peter L. Bartlett, Ling Huang, and Nina Taft. Learning in a large function space: Privacy-preserving mechanisms for SVM learning. *CoRR*, abs/0911.5708, 2009. URL: <http://arxiv.org/abs/0911.5708>.
 - [19] Adam Smith, Jalaj Upadhyay, and Abhradeep Thakurta. Is interaction necessary for distributed private learning? *IEEE Symposium on Security and Privacy*, 2017.
 - [20] Shuang Song, Kamalika Chaudhuri, and Anand D. Sarwate. Stochastic gradient descent with differentially private updates. In *IEEE Global Conference on Signal and Information Processing, GlobalSIP 2013, Austin, TX, USA, December 3-5, 2013*, pages 245–248, 2013. URL: <http://dx.doi.org/10.1109/GlobalSIP.2013.6736861>, doi:10.1109/GlobalSIP.2013.6736861.
 - [21] Oliver Williams and Frank McSherry. Probabilistic inference and differential privacy. In *Advances in Neural Information Processing Systems 23: 24th Annual Conference on Neural Information Processing Systems 2010. Proceedings of a meeting held 6-9 December 2010, Vancouver, British Columbia, Canada.*, pages 2451–2459, 2010. URL: <http://papers.nips.cc/paper/3897-probabilistic-inference-and-differential-privacy>.
 - [22] Steven Wu, Katrina Ligett, Seth Neel, Aaron Roth, and Bo Waggoner. Code for "accuracy first: Selecting a differentially private level for accuracy-constrained erm". [Computer code] 3371774, Zenodo, 2019. URL: <https://doi.org/10.5281/zenodo.3371774>, doi:10.5281/zenodo.3371774.

APPENDIX A. MISSING DETAILS AND PROOFS

A.1. AboveThreshold.

Proof of Lemma 2.8. Let D, D' be neighboring databases. We will instead analyze the algorithm that outputs the entire prefix f_1, \dots, f_t when stopping at time t . Because IAT is a post-processing of this algorithm, and privacy can only be improved under post-processing, this suffices to prove the theorem. We wish to show for all outcomes $o = (t, f_1, \dots, f_t)$:

$$\Pr[\text{IAT}(D) = (t, f_1, f_2, \dots, f_t)] \leq e^{\varepsilon_A + \varepsilon_t} \Pr[\text{IAT}(D') = (t, f_1, f_2, \dots, f_t)].$$

We have directly from the privacy guarantee of InteractiveAboveThreshold that for every *fixed* sequence of queries f_1, \dots, f_t :

$$\Pr[\text{IAT}(D) = t \mid f_1, \dots, f_t] \leq e^{\varepsilon_A} \Pr[\text{IAT}(D') = t \mid f_1, \dots, f_t] \quad (\text{A.1})$$

because the guarantee of InteractiveAboveThreshold is quantified over all data-independent sequences of queries f_1, \dots, f_T , and by definition of the algorithm, the probability of stopping at time t is independent of the identity of any query $f_{t'}$ for $t' > t$.

Now we can write:

$$\Pr[\text{IAT}(D) = t, f_1, \dots, f_t] = \Pr[\text{IAT}(D) = t \mid f_1, \dots, f_t] \Pr[M(D) = f_1, \dots, f_t].$$

By assumption, M is prefix-private, in particular, for fixed t and any f_1, \dots, f_t :

$$\Pr[M(D) = f_1, \dots, f_t] \leq e^{\varepsilon_t} \Pr[M(D') = f_1, \dots, f_t]$$

Thus,

$$\begin{aligned} \frac{\Pr[\text{IAT}(D) = t, f_1, \dots, f_t]}{\Pr[\text{IAT}(D') = t, f_1, \dots, f_t]} &= \frac{\Pr[\text{IAT}(D) = t \mid f_1, \dots, f_t]}{\Pr[\text{IAT}(D') = t \mid f_1, \dots, f_t]} \frac{\Pr[M(D) = f_1, \dots, f_t]}{\Pr[M(D') = f_1, \dots, f_t]} \\ &\leq e^{\varepsilon_A} \cdot e^{\varepsilon_t} = e^{\varepsilon_A + \varepsilon_t}, \end{aligned}$$

as desired. \square

We also include the following utility theorem. We say that an instantiation of InteractiveAboveThreshold is (α, β) accurate with respect to a threshold W and stream of queries f_1, \dots, f_T if except with probability at most γ , the algorithm outputs a query f_t only if $f_t(D) \geq W - \alpha$.

Theorem A.1 . *For any sequence of 1-sensitive queries f_1, \dots, f_T such InteractiveAboveThreshold is (α, β) -accurate for*

$$\alpha = \frac{8\Delta(\log(T) + \log(2/\gamma))}{\varepsilon}.$$

A.2. Doubling Method. We now formally describe the DOUBLINGMETHOD discussed in Section 1 and Section 3, and give a formal ex-post privacy analysis. Let $\theta^* = \operatorname{argmin}_{\theta \in \mathbb{R}^p} L(\theta)$. DOUBLINGMETHOD accepts a list of privacy levels $\varepsilon_1 < \varepsilon_2 < \dots < \varepsilon_T$, where $\varepsilon_i = 2\varepsilon_{i-1}$. We show in Claim B.1 that 2 is the optimal factor to scale ε by. It also takes in a failure probability γ , and a black-box private ERM mechanism M that has the following guarantee: Fixing a dataset D , M takes as input D and a privacy level ε_i , and generates an ε_i -differentially private hypothesis θ_i , such that the query $f^i(D) = L(D, \theta^*) - L(D, \theta_i)$ has ℓ_1 sensitivity at most Δ .

Algorithm 4 Doubling Method: DOUBLINGMETHOD($D, \{\varepsilon_1, \dots, \varepsilon_T\}, M, \alpha, \gamma$)

Input: private dataset D , accuracy α , failure probability γ , mechanism M

for each $t = 1, \dots, T$ **do**
 Generate $\theta_t \leftarrow M(D)_t$
 Let $f^t(D) = L(D, \theta^*) - L(D, \theta_t)$
 Generate $w_t \sim \text{Lap}\left(\frac{\alpha}{2 \log(\frac{T}{\gamma})}\right)$
 if $f^t(D) + w_t \geq -\alpha/2$: **then** Output (t, f^t) ; **Halt.**
Output $T + 1, \theta^*$.

Theorem A.2 . For $k \leq T$, define the privacy loss function $\mathcal{E}(k, \theta_k) = \frac{2k\Delta \log(T/\gamma)}{\alpha} + (2^k - 1)\varepsilon_1$, and $\mathcal{E}(T + 1, \theta^*) = \infty$. Then DOUBLINGMETHOD is \mathcal{E} -ex-post differentially private, and is $1 - \gamma$ accurate.

Proof. Since if the algorithm reaches step $T + 1$ it outputs the true minimizer which has error $0 < \alpha$, it could only fail to output a hypothesis with error less than α if it stops at $i \leq T$. DOUBLINGMETHOD only stops early if the noisy query is greater than $-\alpha/2$; or $f^i(D) + w_i \geq -\alpha/2$. But $f^i(D) \leq -\alpha$, which forces $w_i \geq \alpha/2$. By properties of the Laplace distribution, $\Pr[w_i \geq \alpha/2] = \frac{1}{2} \exp(\frac{-\alpha}{2} \frac{2 \log(\frac{T}{\gamma})}{\alpha}) = \gamma/T$. Hence by union bound over T the total failure probability is at most γ .

By the assumption, generating the k^{th} private hypothesis incurs privacy loss $\varepsilon_1 * 2^{k-1}$. By the Laplace mechanism, evaluating the error of the sensitivity Δ query f^i is $\frac{2\Delta \log(T/\gamma)}{\alpha}$ -differentially private. Theorem 3.6 in (17) then says that the ex-post privacy loss of outputting $k \leq T$ is $\sum_{i=1}^k [\varepsilon_1 * 2^{i-1} + \frac{2\Delta \log(T/\gamma)}{\alpha}] = \frac{2k\Delta \log(T/\gamma)}{\alpha} + (2^k - 1)\varepsilon_1$, as desired. \square

Remark A.3. In practice, the private empirical risk minimization mechanism M may not always output a hypothesis that leads to queries with uniformly bounded ℓ_1 sensitivity. In this case, a projection that scales down, the hypothesis norm can be applied prior to evaluating the private query error. For a discussion of scaling the norm down refer to the experiments section of the appendix.

A.3. Ridge Regression. In this subsection, we let $\ell(\theta, (X_i, y_i)) = \frac{1}{2}(y_i - \langle \theta, X_i \rangle)^2$, and the empirical loss over the data set is defined as

$$L(D, \theta) = \frac{1}{2n} \|y - X\theta\|_2^2 + \frac{\lambda \|\theta\|_2^2}{2},$$

where X denotes the $(n \times p)$ matrix with row vectors X_1, \dots, X_n and $y = (y_1, \dots, y_n)$. We assume that for each i , $\|X_i\|_1 \leq 1$ and $|y_i| \leq 1$. For simplicity, we will sometimes write $L(\theta)$ for $L(D, \theta)$.

First, we show that the unconstrained optimal solution in ridge regression has bounded norm.

Lemma A.4. Let $\theta^* = \arg\min_{\theta \in \mathbb{R}^d} L(\theta)$. Then $\|\theta^*\|_2 \leq \frac{1}{\sqrt{\lambda}}$.

Proof. For any $\theta \in \mathbb{R}^p$, $L(\theta^*) \leq L(\theta)$. In particular for $\theta = \mathbf{0}$,

$$L(\theta^*) \leq L(\mathbf{0}) = \sum_{i=1}^n \frac{1}{2n} \ell((X_i, y_i), 0) \leq \frac{1}{2}.$$

Note that for any θ , $\ell((X_i, y_i), \theta) \geq 0$, so this means $L(\theta^*) \geq \frac{\lambda}{2} \|\theta^*\|_2^2$, which forces $\frac{\lambda}{2} \|\theta^*\|_2^2 \leq \frac{1}{2}$, and so $\|\theta^*\|_2 \leq \frac{1}{\sqrt{\lambda}}$ as desired. \square

The following claim provides a bound on the sensitivity for the excess risk, which are the queries we send to `InteractiveAboveThreshold`.

Claim A.5. *Let C be a bounded convex set in \mathbb{R}^p with $\|C\|_2 \leq M$. Let D and D' be a pair of adjacent datasets, and let $\theta^* = \operatorname{argmin}_{\theta \in C} L(\theta, D)$ and $\theta^\bullet = \operatorname{argmin}_{\theta \in C} L(\theta, D')$. Then for any $\theta \in C$,*

$$|(L(\theta, D) - L(\theta^*, D)) - (L(\theta, D') - L(\theta^\bullet, D'))| \leq \frac{(M+1)^2}{n}.$$

The following lemma provides a bound on the ℓ_1 sensitivity for the matrix $X^\top X$ and vector $X^\top y$.

Lemma A.6. *Fix any $i \in [n]$. Let X and Z be two $n \times p$ matrices such that for all rows $j \neq i$, $X_j = Z_j$. Let $y, y' \in \mathbb{R}^n$ such that $y_j = y'_j$ for all $j \neq i$. Then*

$$\|X^\top X - Z^\top Z\|_1 \leq 2 \quad \text{and} \quad \|X^\top y - Z^\top y'\|_1 \leq 2,$$

as long as $\|X_i\|, \|Z_i\|, |y_i|, |y'_i| \leq 1$.

Proof. We can write

$$\begin{aligned} \|X^\top X - Z^\top Z\|_1 &= \left\| \sum_j (X_j^\top X_j - Z_j^\top Z_j) \right\|_1 \\ &= \|X_i^\top X_i - Z_i^\top Z_i\|_1 \\ &\leq \|X_i^\top X_i\|_1 + \|Z_i^\top Z_i\|_1 \\ &= \|X_i\|_1^2 + \|Z_i\|_1^2 \leq 2. \end{aligned}$$

Similarly,

$$\begin{aligned} \|X^\top y - Z^\top y'\|_1 &= \left\| \sum_j (y_j X_j - y'_j Z_j) \right\|_1 \\ &= \|y_i X_i - y'_i Z_i\|_1 \\ &= \|y_i X_i\|_1 + \|y'_i Z_i\|_1 \\ &= \|X_i\|_1 + \|Z_i\|_1 \leq 2. \end{aligned}$$

This completes the proof. \square

Before we proceed to give a formal proof for Theorem 3.1, we will also give the following basic fact about Laplace random vectors.

Claim A.7. *Let $\nu = (\nu_1, \dots, \nu_k) \in \mathbb{R}^k$ such that each ν_i is an independent random variable drawn from the Laplace distribution $\text{Lap}(r)$. Then $\mathbb{E}[\|\nu\|_2] \leq \sqrt{2kr}$.*

Proof. By Jensen's inequality,

$$\mathbb{E}[\|\nu\|_2] = \mathbb{E}\left[\sqrt{\sum_i \nu_i^2}\right] \leq \sqrt{\mathbb{E}\left[\sum_i \nu_i^2\right]}.$$

Note that by linearity of expectation and the variance of the Laplace distribution

$$\mathbb{E}\left[\sum_i \nu_i^2\right] = \sum_i \mathbb{E}[\nu_i^2] = \sum_i 2r^2 = 2kr^2.$$

Therefore, we have $\mathbb{E}[\|\nu\|_2] \leq \sqrt{2kr}$. \square

Proof of Theorem 3.1. In the algorithm, we compute $Z = X^\top X + B$ and $z = X^\top y + b$, where the entries of B and b are drawn i.i.d. from $\text{Lap}(4/\varepsilon)$. Note that the output θ_p is simply a post-processing of the noisy matrix Z and vector z . Furthermore, by Lemma A.6, the joint vector (Z, z) is has sensitivity bounded by 4 with respect to ℓ_1 norm. Therefore, the mechanism satisfies ε -differential privacy by the privacy guarantee of the Laplace mechanism.

Let $M = \sqrt{1/\lambda}$ and $L_p(\theta) = \frac{1}{2n}(-2\langle z, \theta \rangle) + \frac{1}{2n}(\theta^\top Z \theta) + \frac{\lambda \|\theta\|_2^2}{2}$. Observe that $\theta_p = \text{argmin}_{\theta \in C} L_p(\theta)$. Our goal is to bound $L(\theta_p) - L(\theta^*)$, which can be written as follows

$$\begin{aligned} L(\theta_p) - L(\theta^*) &= L(\theta_p) - L_p(\theta_p) + L_p(\theta_p) - L_p(\theta^*) + L_p(\theta^*) - L(\theta^*) \\ &\leq L(\theta_p) - L_p(\theta_p) + L_p(\theta^*) - L(\theta^*) \\ &= \frac{1}{2n}(2\langle b, \theta_p \rangle - \theta_p^\top B \theta_p) - \frac{1}{2n}(2\langle b, \theta^* \rangle - (\theta^*)^\top B \theta^*) \end{aligned}$$

Moreover, $\langle b, \theta_p \rangle \leq \|b\|_2 \|\theta_p\|_2 \leq M \|b\|_2$ and

$$\begin{aligned} -\theta_p^\top B \theta_p &= -\sum_{(s,t) \in [p]^2} B_{st}(\theta_p)_s(\theta_p)_t \\ &\leq \left(\sum_{(s,t)} B_{st}^2\right)^{1/2} \left(\sum_{s,t} (\theta_p)_s^2 (\theta_p)_t^2\right)^{1/2} \\ &= \|B\|_F \left[\left(\sum_s (\theta_p)_s^2\right)^2\right]^{1/2} \\ &\leq \|B\|_F M^2 \end{aligned}$$

By Claim A.7, we also have $\mathbb{E} [\|B\|_F] \leq 4\sqrt{2}p/\varepsilon$ and $\mathbb{E} [\|b\|_2] \leq 4\sqrt{2}p/\varepsilon$. Finally,

$$\begin{aligned} \mathbb{E} [L(\theta_p) - L(\theta^*)] &\leq \mathbb{E} \left[\frac{1}{2n} (2\langle b, \theta_p \rangle - \theta_p^\top B \theta_p) - \frac{1}{2n} (2\langle b, \theta^* \rangle - (\theta^*)^\top B \theta^*) \right] \\ &= \mathbb{E} \left[\frac{2\langle b, \theta_p \rangle - \theta_p^\top B \theta_p}{2n} \right] \\ &\leq \frac{\mathbb{E} [2M\|b\|_2] + \mathbb{E} [M^2\|B\|_F]}{2n} \\ &\leq \frac{4\sqrt{2}(2\sqrt{p}M + pM^2)}{n\varepsilon} \end{aligned}$$

which recovers our stated bound. \square

Next, we will also provide a theoretical result for applying output perturbation (with Laplace noise) to the ridge regression problem. This will provides us the “theory curve” for output perturbation in ridge regression plot of Figure 1a.

First, the following sensitivity bound on the optimal solution for L follows directly from the strong convexity of L .

Lemma A.8. *Let C be a bounded convex set in \mathbb{R}^p with $\|C\|_2 \leq M$. Let D and D' be a pair of neighboring datasets, and let $\theta^* = \operatorname{argmin}_{\theta \in C} L(\theta, D)$ and $\theta^\bullet = \operatorname{argmin}_{\theta \in C} L(\theta, D')$. Then $\|\theta^* - \theta^\bullet\|_1 \leq (M + 1)\sqrt{\frac{p}{n\lambda}}$.*

Theorem A.9 . *Let $\varepsilon > 0$ and C be a bounded convex set with $\|C\|_2 \leq \sqrt{1/\lambda}$. Let $r = (\sqrt{1/\lambda} + 1)\sqrt{p/(n\lambda)}/\varepsilon$. Consider the following mechanism \mathcal{M} that for any input dataset D first computes the optimal solution $\theta^* = \operatorname{argmin}_{\theta \in C} L(\theta)$, and then outputs $\theta_p = \theta^* + b$, where b is a random vector with its entries drawn i.i.d. from $\operatorname{Lap}(r)$. Then \mathcal{M} satisfies ε -differential privacy, and θ_p satisfies*

$$\mathbb{E}_b [L(\theta_p) - L(\theta^*)] \leq \left(\frac{1}{n} + \lambda \right) \frac{(\sqrt{1/\lambda} + 1)^2 p^2}{n\lambda\varepsilon^2}.$$

Proof. The privacy guarantee follows directly from the use of Laplace mechanism and the ℓ_1 sensitivity bound in Lemma A.8.

For each data point $d_i = (X_i, y_i)$, we have

$$\begin{aligned} (y_i - \langle \theta_p, X_i \rangle)^2 - (y_i - \langle \theta^*, X_i \rangle)^2 &= (\langle \theta_p, X_i \rangle)^2 - (\langle \theta^*, X_i \rangle)^2 - 2\langle b, X_i \rangle \\ &= b^\top (X_i^\top X_i) b + (\theta^*)^\top (X_i^\top X_i) b + b^\top (X_i^\top X_i) \theta^* - 2\langle b, X_i \rangle \end{aligned}$$

Since each entry in b has mean 0, we can simplify the expectation as

$$\begin{aligned} \mathbb{E} [(y_i - \langle \theta_p, X_i \rangle)^2 - (y_i - \langle \theta^*, X_i \rangle)^2] &= \mathbb{E} [b^\top (X_i^\top X_i) b] \\ &= \mathbb{E} [\langle b, X_i \rangle^2] \\ &\leq \mathbb{E} [\|b\|_2^2 \|X_i\|_2^2] \\ &= \mathbb{E} [\|b\|_2^2] \mathbb{E} [\|X_i\|_2^2] \\ &\leq \mathbb{E} [\|b\|_2^2] \leq 2pr^2 \end{aligned}$$

In the following, let $M = \sqrt{1/\lambda}$. We can then bound

$$\begin{aligned}\|\theta_p\|_2^2 - \|\theta^*\|_2^2 &= \sum_{s \in [p]} [(\theta_s + b_s)^2 - \theta_s^2] \\ &= \sum_{s \in [p]} [2\theta_s b_s + b_s^2],\end{aligned}$$

Again, since each b_s is drawn from $\text{Lap}(r)$, we get

$$\begin{aligned}\mathbb{E} [\|\theta_p\|_2^2 - \|\theta^*\|_2^2] &= \mathbb{E} \left[\sum_s b_s^2 \right] \\ &= \sum_s \mathbb{E} [b_s^2] = 2pr^2.\end{aligned}$$

To put all the pieces together and plugging in the value of r , we get

$$\begin{aligned}\mathbb{E}_b [L(\theta_p) - L(\theta^*)] &\leq \left(\frac{1}{2n} + \frac{\lambda}{2} \right) 2pr^2 \\ &= \left(\frac{1}{n} + \lambda \right) \frac{(M+1)^2 p^2}{n\lambda\epsilon^2}\end{aligned}$$

which recovers our stated bound. \square

A.4. Logistic Regression. In this subsection, the input data D consists of n labelled examples $(X_1, y_1), \dots, (X_n, y_n)$, such that for each i , $x_i \in \mathbb{R}^p$, $\|x_i\|_1 \leq 1$, and $y_i \in \{-1, 1\}$.

We consider the logistic loss function: $\ell(\theta, (X_i, y_i)) = \log(1 + \exp(-y_i \theta^\top X_i))$, and our empirical loss is defined as

$$L(\theta, D) = \frac{1}{n} \sum_{i=1}^n \log(1 + \exp(-y_i \theta^\top X_i)) + \frac{\lambda \|\theta\|_2^2}{2}.$$

In output perturbation, the noise needs to scale with the ℓ_1 -sensitivity of the optimal solution, which is given by the following lemma.

Lemma A.10. *Let D and D' be a pair of neighboring datasets. Let $\theta = \arg\min_{w \in \mathbb{R}^p} L(w, D)$ and $\theta' = \arg\min_{w' \in \mathbb{R}^p} L(w', D')$. Then $\|\theta - \theta'\|_1 \leq \frac{2\sqrt{p}}{n\lambda}$.*

Proof of Lemma A.10. By Corollary 8 of (5), we can bound

$$\|\theta - \theta'\|_2 \leq \frac{2}{n\lambda}$$

By the fact that $\|a\|_1 \leq \sqrt{p}\|a\|_2$ for any $a \in \mathbb{R}^p$, we recover the stated result. \square

We will show that the optimal solution for the unconstrained problem has ℓ_2 norm no more than $\sqrt{2\log 2/\lambda}$.

Claim A.11. *The (unconstrained) optimal solution θ^* has norm $\|\theta^*\|_2 \leq \sqrt{\frac{2\log 2}{\lambda}}$.*

Proof. Note that the weight vector $\theta = \vec{0}$ has loss $\log 2$. Therefore, $L(\theta^*) \leq \log 2$. Since the logistic loss is positive, we know that the regularization term

$$\frac{\lambda}{2} \|\theta^*\|_2^2 \leq \log 2.$$

It follows that $\|\theta^*\|_2 \leq \sqrt{\frac{2\log 2}{\lambda}}$. □

We will focus on generating hypotheses θ within the set $C = \{a \in \mathbb{R}^p \mid \|a\|_2 \leq \sqrt{2\log 2/\lambda}\}$. Then we can bound the ℓ_1 sensitivity of the excess risk using the following result.

Claim A.12. *Let D and D' be a pair of neighboring datasets. Then for any $\theta \in \mathbb{R}^p$ such that $\|\theta\|_2 \leq M$,*

$$|L(\theta, D) - L(\theta, D')| \leq \frac{2}{n} \log \left(\frac{1 + \exp(M)}{1 + \exp(-M)} \right)$$

The following fact is useful for our utility analysis for the output perturbation method.

Claim A.13. *Fix any data point (x, y) such that $\|x\|_1 \leq 1$ and $y \in \{-1, 1\}$. The logistic loss function $\ell(\theta, (x, y))$ is a 1-Lipschitz function in θ .*

Proof of Theorem 3.3. The privacy guarantee follows directly from the use of Laplace mechanism and the ℓ_1 -sensitivity bound in Lemma A.10. Since the logistic loss function is 1-Lipschitz. For any (x, y) in our domain,

$$|\ell(\theta^*, (x, y)) - \ell(\theta^p, (x, y))| \leq \|\theta^* - \theta^p\|_2 = \|b\|_2.$$

Furthermore,

$$\|\theta_p\|_2^2 - \|\theta^*\|_2^2 = \|\theta^* + b\|_2^2 - \|\theta^*\|_2^2 = 2\langle b, \theta^* \rangle + \|b\|_2^2$$

By Claim A.7 and the property of the Laplace distribution, we know that

$$\mathbb{E}[\|b\|_2] \leq \sqrt{2pr} \quad \text{and} \quad \mathbb{E}[\|b\|_2^2] = 2pr^2.$$

It follows that

$$\begin{aligned} \mathbb{E}_b[L(\theta_p) - L(\theta^*)] &\leq \mathbb{E}_b[\|b\|_2] + \frac{\lambda}{2} \mathbb{E}[\|b\|_2^2] \\ &\leq \sqrt{2pr} + p\lambda r^2 = \frac{2\sqrt{2pr}}{n\lambda\varepsilon} + \frac{4p^2}{n^2\lambda\varepsilon^2}, \end{aligned}$$

which recovers the stated bound. □

Algorithm 5 Output Perturbation with Noise-Reduction:
 OUTPUTNR($D, \{\varepsilon_1, \dots, \varepsilon_T\}, \alpha, \gamma$)

Input: private data set $D = (X, y)$, accuracy parameter α , privacy levels $\varepsilon_1 < \varepsilon_2 < \dots < \varepsilon_T$, and failure probability γ

Let $M = \sqrt{2 \log 2 / \lambda}$

Instantiate Interactive AboveThreshold: $\mathcal{A} = (D, \varepsilon_0, \alpha/2, 2 \log(1 + \exp(M))/(1 + \exp(-M))/(n), \cdot)$ with $\varepsilon_0 = 16\Delta(\log(2T/\gamma))/\alpha$ and $\Delta = 2 \log(1 + \exp(M))/(1 + \exp(-M))/(n)$

Let $C = \{a \in \mathbb{R}^p \mid \|a\|_2 \leq \sqrt{1/\lambda}\}$ and $\theta^* = \operatorname{argmin}_{\theta \in \mathbb{R}^p} L(\theta)$

Generate hypotheses: $\{\theta^t\} = \text{NR}(\theta^*, \frac{2\sqrt{p}}{n\lambda}, \{\varepsilon_1, \dots, \varepsilon_T\})$

for $t = 1, \dots, T$: **do**

if $\|\theta^t\|_2 \leq M$ **then** Set $\theta^t = M(\theta^t / \|\theta^t\|_2)$

 Let $f^t(D) = L(D, \theta^*) - L(D, \theta^t)$

 Query \mathcal{A} with f^t

if yes **then** **Output** (t, θ^t)

Output: (\perp, θ^*)

We include the full details of OUTPUTNR in Algorithm 5.

APPENDIX B. EXPERIMENTS

B.1. Parameters and data. For simplicity and to avoid over-fitting, we fixed the following parameters for both experiments:

- $n = 100,000$ (number of data points)
- $\lambda = 0.005$ (regularization parameter)
- $\gamma = 0.10$ (requested failure probability)
- $\varepsilon_1 = 4E$, where E is the inversion of the theory guarantee for the underlying algorithm. For example in the logistic regression setting where the algorithm is Output Perturbation, E is the value such that setting $\varepsilon = E$ guarantees **expected** excess risk of at most α .
- $\varepsilon_T = 1.0/n$.
- $\alpha = 0.005, 0.010, 0.015, \dots, 0.200$ (requested excess error bound).

For NoiseReduction, we choose $T = 1000$ (maximum number of iterations) and set $\varepsilon_t = \varepsilon_1 r^t$ for the appropriate r , i.e. $r = \left(\frac{\varepsilon_T}{\varepsilon_1}\right)^{1/T}$.

For the Doubling method, T is equal to the number of doubling steps until ε_t exceeds ε_T , i.e. $T = \lceil \log_2(\varepsilon_1/\varepsilon_T) \rceil$.

Features, labels, and transformations. The Twitter dataset has $p = 77$ features (dimension of each x), relating to measurements of activity relating to a posting; the label y is a measurement of the “buzz” or success of the posting. Because general experience suggests that such numbers likely follow a heavy-tailed distribution, we transformed the labels by $y \mapsto \log(1 + y)$ and set the tasks of predicting the transformed label.

The KDD-99 Cup dataset has $p = 38$ features relating to attributes of a network connection such as duration of connection, number of bytes sent in each direction, binary attributes, etc. The goal is to classify connections as innocent or malicious, with malicious connections broken down into further subcategories. We transformed three attributes

containing likely heavy-tailed data (the first three mentioned above) by $x_i \mapsto \log(1 + x_i)$, dropped three columns containing textual categorical data, and transformed the labels into 1 for any kind of malicious connection and 0 for an innocent one. (The feature length $p = 38$ is after dropping the text columns.)

For both datasets, we transformed the data by renormalizing to maximum $L1$ -norm 1. That is, we computed $M = \max_i \|x_i\|_1$, and transformed each $x_i \mapsto x_i/M$. In the case of the Twitter dataset, we did the same (separately) for the y labels. This is *not* a private operation (unlike the previous ones) on the data, as it depends precisely on the maximum norm. We do not consider the problem of privately ensuring bounded-norm data, as it is orthogonal to the questions we study.

The code for the experiments is implemented in python3 using the numpy and scikit-learn libraries.

B.2. Additional results. Figure 2 plots the empirical accuracies of the output hypotheses, to ensure that the algorithms are achieving their theoretical guarantees. In fact, they do significantly better, which is reasonable considering the private testing methodology: set a threshold significantly below the goal α , add independent noise to each query, and accept only if the query plus noise is smaller than the threshold. Combined with the requirement to use tail bounds, the accuracies tend to be significantly smaller than α and with significantly higher probability than $1 - \gamma$. (Recall: this is not necessarily a good thing, as it probably costs a significant amount of extra privacy.)

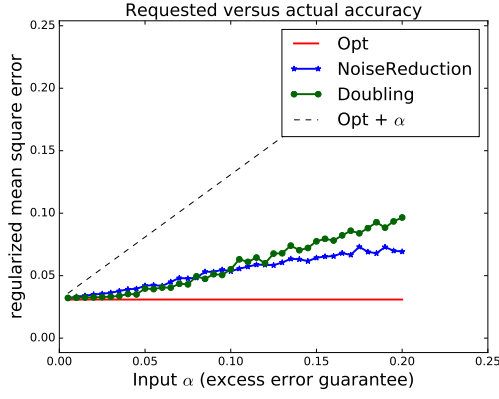
Figure 3 shows the breakdown in privacy losses between the “privacy test” and the “hypothesis generator”. In the case of NoiseReduction, these are AboveThreshold’s ε_A and the ε_t of the private method, Covariance Perturbation or Output Perturbation. In the case of Doubling, these are the accrued ε due to tests at each step and due to Covariance Perturbation or Output Perturbation for outputting the hypotheses.

This shows the majority of the privacy loss is due to testing for privacy levels. One reason why might be that the cost of privacy tests depends heavily on certain constants, such as the norm of the hypothesis being tested. This norm is upper-bounded by a theoretical maximum which is used, but a smaller maximum would allow for significantly higher computed privacy levels for the same algorithm. In other words, the analysis might be loose compared to an analysis that knows the norms of the hypotheses, although this is a private quantity. Figure 4 supports the conclusion that generally, the theoretical maximum was very pessimistic in our cases. Note that a tenfold reduction in norm gives a tenfold reduction in privacy level for logistic regression, where sensitivity is linear in maximum norm; and a *hundred-fold* reduction for ridge regression.

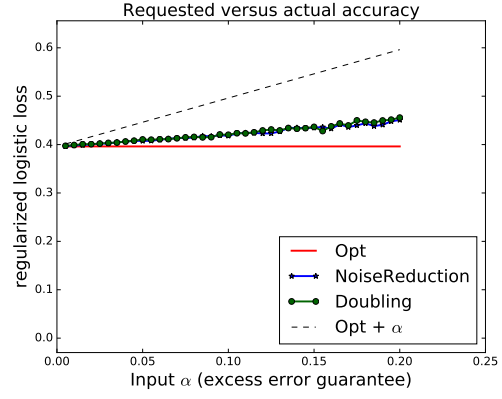
B.3. Supporting theory.

Claim B.1. *For the “doubling method”, the factor 2 increase in ε at each time step gives the optimal worst case ex post privacy loss guarantee.*

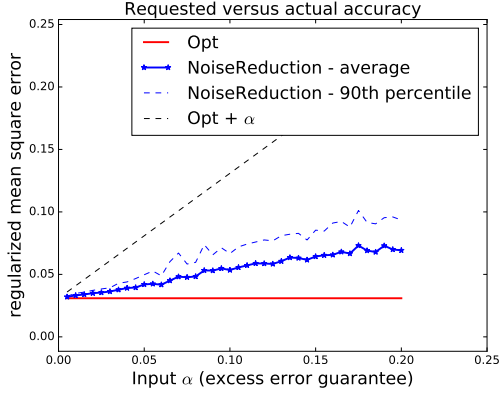
Proof. In a given setting, suppose ε^* is the “final” level of privacy at which the algorithm would halt. With a factor $1/r$ increase for $r < 1$, the final loss may be as large as ε^*/r . The



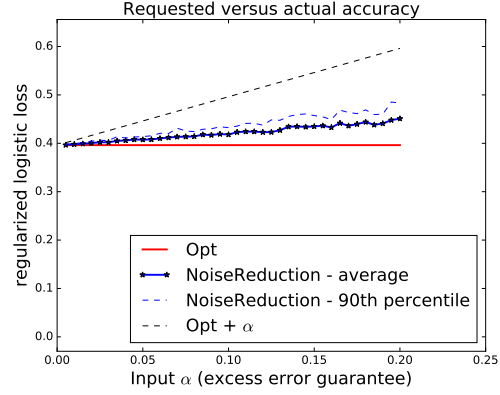
(A) Linear (ridge) regression.



(B) Regularized logistic regression.



(C) Linear (ridge) regression.



(D) Regularized logistic regression.

FIGURE 2. **Empirical accuracies.** The dashed line shows the requested accuracy level, while the others plot the actual accuracy achieved. Due most likely due to a pessimistic analysis and the need to set a small testing threshold, accuracies are significantly better than requested for both methods.

total loss is the sum of that loss and all previous losses, i.e. if t steps were taken:

$$\begin{aligned}
 (\varepsilon^*/r) + r \cdot (\varepsilon^*/r) + \cdots + r^{t-1} \cdot (\varepsilon^*/r) &= (\varepsilon^*/r) \sum_{j=0}^{t-1} r^j \\
 &\rightarrow (\varepsilon^*/r) \sum_{j=0}^{\infty} r^j \\
 &= \frac{\varepsilon^*}{r(1-r)} \\
 &\geq 4\varepsilon^*.
 \end{aligned}$$

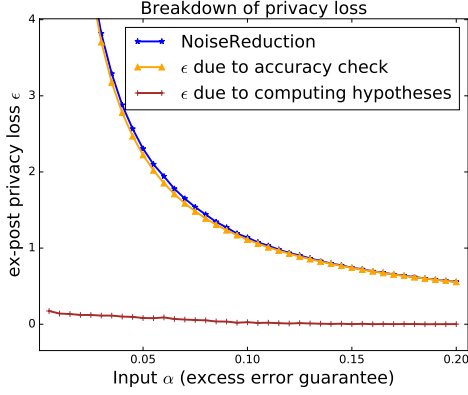
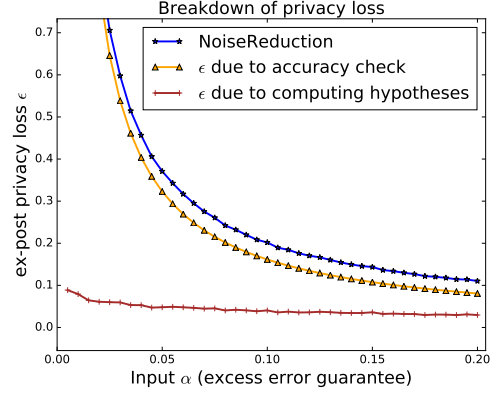
(A) **Linear (ridge) regression.**(B) **Regularized logistic regression.**

FIGURE 3. **Privacy breakdowns.** Shows the amount of empirical privacy loss due to computing the hypotheses themselves and the losses due to testing their accuracies.

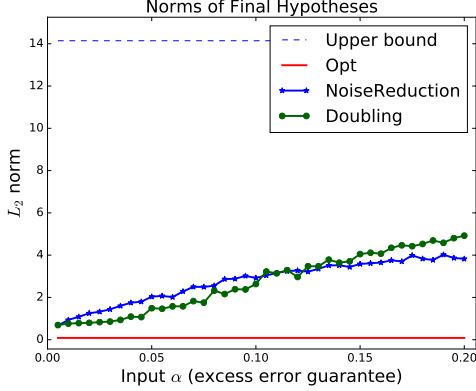
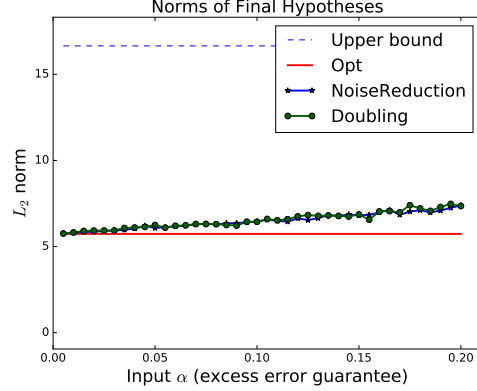
(A) **Linear (ridge) regression.**(B) **Regularized logistic regression.**

FIGURE 4. **L_2 norms of final hypotheses.** Shows the average L_2 norm of the output $\hat{\theta}$ for each method, versus the theoretical maximum of $1/\sqrt{\lambda}$ in the case of ridge regression and $\sqrt{2\log(2)/\lambda}$ in the case of regularized logistic regression.

The final inequality implies that setting $r = 0.5$ and $(1/r) = 2$ is optimal. The asymptotic \rightarrow is justified by noting that the starting ε_1 may be chosen arbitrarily small, so there exist parameters that exceed the value of that summation for any finite t ; and the summation limits to $\frac{1}{1-r}$ as $t \rightarrow \infty$. \square