

Sample Complexity Bounds on Differentially Private Learning via Communication Complexity

Vitaly Feldman* David Xiao†

September 16, 2015

Abstract

In this work we analyze the sample complexity of classification by differentially private algorithms. Differential privacy is a strong and well-studied notion of privacy introduced by [Dwork et al. \[2006\]](#) that ensures that the output of an algorithm leaks little information about the data point provided by any of the participating individuals. Sample complexity of private PAC and agnostic learning was studied in a number of prior works starting with [\[Kasiviswanathan et al., 2011\]](#). However, a number of basic questions still remain open [\[Beimel et al., 2010, Chaudhuri and Hsu, 2011, Beimel et al., 2013a,b\]](#), most notably whether learning with privacy requires more samples than learning without privacy.

We show that the sample complexity of learning with (pure) differential privacy can be arbitrarily higher than the sample complexity of learning without the privacy constraint or the sample complexity of learning with approximate differential privacy. Our second contribution and the main tool is an equivalence between the sample complexity of (pure) differentially private learning of a concept class C (or $\text{SCDP}(C)$) and the randomized one-way communication complexity of the evaluation problem for concepts from C . Using this equivalence we prove the following bounds:

- $\text{SCDP}(C) = \Omega(\text{LDim}(C))$, where $\text{LDim}(C)$ is the Littlestone’s dimension characterizing the number of mistakes in the online-mistake-bound learning model [\[Littlestone, 1987\]](#). Known bounds on $\text{LDim}(C)$ then imply that $\text{SCDP}(C)$ can be much higher than the VC-dimension of C .
- For any t , there exists a class C such that $\text{LDim}(C) = 2$ but $\text{SCDP}(C) \geq t$.
- For any t , there exists a class C such that the sample complexity of (pure) α -differentially private PAC learning is $\Omega(t/\alpha)$ but the sample complexity of the approximate (α, β) -differentially private PAC learning is $O(\log(1/\beta)/\alpha)$. This resolves an open problem from [\[Beimel et al., 2013b\]](#).

1 Introduction

In machine learning tasks, the training data often consists of information collected from individuals. This data can be highly sensitive, for example in the case of medical or financial information, and therefore privacy-preserving data analysis is becoming an increasingly important area of study in machine learning, data mining and statistics [\[Dwork and Smith, 2009, Sarwate and Chaudhuri, 2013, Dwork and Roth, 2014\]](#).

In this work we focus on the task of learning to classify from labeled examples. Two standard and closely related models of this task are PAC learning [\[Valiant, 1984\]](#) and agnostic [\[Haussler, 1992, Kearns et al., 1994\]](#) learning. In the PAC learning model the algorithm is given random examples in which each point is sampled

*IBM Research - Almaden. E-mail: vitaly@post.harvard.edu. Part of this work done while visiting LIAFA, Université Paris 7.

†CNRS, Université Paris 7. E-mail: dxiao@liafa.univ-paris-diderot.fr. Part of this work done while visiting Harvard’s Center for Research on Computation and Society (CRCS).

⁰Preliminary version of this work has appeared in Conference on Learning Theory (COLT), 2014

i.i.d. from some unknown distribution over the domain and is labeled by an unknown function from a set of functions C (called concept class). In the agnostic learning model the algorithm is given examples sampled i.i.d. from an arbitrary (and unknown) distribution over labeled points. The goal of the learning algorithm in both models is to output a hypothesis whose prediction error on the distribution from which examples are sampled is not higher (up to an additive ε) than the prediction error of the best function in C (which is 0 in the PAC model). See [Section 2.1](#) for formal definitions.

We rely on the well-studied differential privacy model of privacy. Differential privacy gives a formal semantic guarantee of privacy, saying intuitively that no single individual’s data has too large of an effect on the output of the algorithm, and therefore observing the output of the algorithm does not leak much information about an individual’s private data [[Dwork et al., 2006](#)] (see [Section 2.2](#) for the formal definition). The downside of this desirable guarantee is that for some problems achieving it has an additional cost: both in terms of the number of examples, or sample complexity, and computation.

The cost of differential privacy in PAC and agnostic learning was first studied by [Kasiviswanathan et al. \[2011\]](#). They showed that the sample complexity¹ of differentially privately learning a concept class C over domain X , denoted by $\text{SCDP}(C)$, is $O(\log(|C|))$ and left open the natural question of whether $\text{SCDP}(C)$ is different from the VC dimension of C which, famously, characterizes the sample complexity of learning C (without privacy constraints). By Sauer’s lemma, $\log(|C|) = O(\text{VC}(C) \cdot \log(|X|))$ and therefore the multiplicative gap between these two measures can be as large as $\log(|X|)$.

Subsequently, [Beimel et al. \[2010\]](#) showed that there exists a large concept class, specifically single points, for which the sample complexity of learning with privacy is a constant. They also show that differentially private *proper* learning (the output hypothesis has to be from C) of single points Point_b and threshold functions Thr_b on the set $I_b = \{0, 1, \dots, 2^b - 1\}$ requires $\Omega(b)$ samples. These results demonstrate that the sample complexity can be lower than $O(\log(|C|))$ and also that lower bounds on the sample complexity of proper learning do not necessarily apply to non-proper learning that we consider here. A similar lower bound on proper learning of thresholds on an interval was given by [Chaudhuri and Hsu \[2011\]](#) in a continuous setting where the sample complexity becomes infinite. They also showed that the sample complexity can be reduced to essentially $\text{VC}(C)$ by either adding distributional assumptions or by requiring only the privacy of the labels.

The upper bound of [Beimel et al. \[2010\]](#) is based on an observation from [[Kasiviswanathan et al., 2011](#)] that if there exists a class of functions H such that for every $f \in C$ and every distribution \mathcal{D} over the domain, there exists $h \in H$ such that $\Pr_{x \sim \mathcal{D}}[f(x) \neq h(x)] \leq \varepsilon$ then the sample complexity of differentially private PAC learning with error 2ε can be reduced to $O(\log(|H|)/\varepsilon)$. They refer to such H as an ε -representation of C , and define the (deterministic) ε -representation dimension of C , denoted as $\text{DRDim}_\varepsilon(C)$, as $\log(|H|)$ for the smallest H that ε -represents C . We note that this natural notion can be seen as a distribution-independent version of the standard notion of ε -covering of C in which the distribution over the domain is fixed [[e.g. Benedek and Itai, 1991](#)].

[Beimel et al. \[2013a\]](#) then defined a probabilistic relaxation of ε -representation defined as follows. A distribution \mathcal{H} over sets of boolean functions on X is said to (ε, δ) -probabilistically represent C if for every $f \in C$ and distribution \mathcal{D} over X , with probability $1 - \delta$ over the choice of $H \stackrel{R}{\leftarrow} \mathcal{H}$, there exists $h \in H$ such that $\Pr_{x \sim \mathcal{D}}[h(x) \neq f(x)] \leq \varepsilon$. The (ε, δ) -probabilistic representation dimension $\text{PRDim}_{\varepsilon, \delta}(C)$ is the minimal $\max_{H \in \text{supp}(\mathcal{H})} \log |H|$, where the minimum is over all \mathcal{H} that (ε, δ) -probabilistically represent C . [Beimel et al. \[2013a\]](#) demonstrated that $\text{PRDim}_{\varepsilon, \delta}(C)$ characterizes the sample complexity of differentially private PAC learning. In addition, they show that DRDim can upper-bounded by PRDim as $\text{DRDim}(C) = O(\text{PRDim}(C) + \log \log(|X|))$, where we omit ε and δ when they are equal to $1/4$.

[Beimel et al. \[2013b\]](#) consider PAC learning with *approximate* (α, β) -differential privacy where the privacy guarantee holds with probability $1 - \beta$ (the basic notion is also referred to as *pure* to distinguish it from the approximate version). They show that Thr_b can be PAC learned using $O(16^{\log^*(b)} \cdot \log(1/\beta))$ samples

¹For now we ignore the dependence on other parameters and consider them to be small constants.

(α is a constant as before). Their algorithm is proper so this separates the sample complexity of pure differentially private proper PAC learning from the approximate version. This work leaves open the question of whether such a separation can be proved for (non-proper) PAC learning.

1.1 Our results

In this paper we resolve the open problems described above. In the process we also establish a new relation between SCDP and Littlestone’s dimension, a well-studied measure of sample complexity of online learning [Littlestone, 1987] (see Section 2.5 for the definition). The main ingredient of our work is a characterization of DRDim and PRDim in terms of randomized one-way communication complexity of associated evaluation problems [Kremer et al., 1999]. In such a problem Alice is given as input a function $f \in C$ and Bob is given an input $x \in X$. Alice sends a single message to Bob, and Bob’s goal is to compute $f(x)$. The question is how many bits Alice must communicate to Bob in order for Bob to be able to compute $f(x)$ correctly, with probability at least $2/3$ over the randomness used by Alice and Bob.

In the standard or “private-coin” version of this model, Alice and Bob each have their own source of random coins. The minimal number of bits needed to solve the problem for all $f \in C$ and $x \in X$ is denoted by $R^{\rightarrow}(C)$. In the stronger “public coin” version of the model, Alice and Bob share the access to the same source of random coins. The minimal number of bits needed to evaluate C (with probability at least $2/3$) in this setting is denoted by $R^{\rightarrow, \text{pub}}(C)$. See Section 2.4 for formal definitions.

We show that these communication problems are equivalent to deterministic and probabilistic representation dimensions of C and, in particular, $\text{SCDP}(C) = \theta(R^{\rightarrow, \text{pub}}(C))$ (for clarity we omit the accuracy and confidence parameters, see Theorem 3.1 and Theorem 3.2 for details).

Theorem 1.1. $\text{DRDim}(C) = \Theta(R^{\rightarrow}(C))$ and $\text{PRDim}(C) = \Theta(R^{\rightarrow, \text{pub}}(C))$.

The evaluation of threshold functions on a (discretized) interval I_b corresponds to the well-studied “greater than” function in communication complexity denoted as GT. $\text{GT}_b(x, y) = 1$ if and only if $x > y$, where $x, y \in \{0, 1\}^b$ are viewed as binary representations of integers. It is known that $R^{\rightarrow, \text{pub}}(\text{GT}_b) = \Omega(b)$ [Miltersen et al., 1998]. By combining this lower bound with Theorem 1.1 we obtain a class whose VC dimension is 1 yet it requires at least $\Omega(b)$ samples to PAC learn differentially privately.

This equivalence also shows that some of the known results in [Beimel et al., 2010, 2013a] are implied by well-known results from communication complexity, sometimes also giving simpler proofs. For example (1) the constant upper bound on the sample complexity of single points follows from the communication complexity of the equality function and (2) the bound $\text{DRDim}(C) = O(\text{PRDim}(C) + \log \log(|X|))$ follows from the classical result of Newman [1991] on the relationship between the public and private coin models. See Section 3.1 for more details and additional examples.

Our second contribution is a relationship of $\text{SCDP}(C)$ (via the equivalences with $R^{\rightarrow, \text{pub}}(C)$) to Littlestone’s [1987] dimension of C . Specifically, we prove

Theorem 1.2. 1. $R^{\rightarrow, \text{pub}}(C) = \Omega(\text{LDim}(C))$.

2. For any t , there exists a class C such that $\text{LDim}(C) = 2$ but $R^{\rightarrow, \text{pub}}(C) \geq t$.

The first result follows from a natural reduction to the augmented index problem, which is well-studied in communication complexity [Bar-Yossef et al., 2004]. While new in our context, the relationship of Littlestone’s dimension to quantum communication complexity was shown by Zhang [2011]. Together with numerous known bounds on LDIM [e.g. Littlestone, 1987, Maass and Turán, 1994b], our result immediately yields a number of new lower bounds on SCDP. In particular, results of Maass and Turán [1994b] imply that linear threshold functions over I_b^d require $\Omega(d^2 \cdot b)$ samples to learn differentially privately. This implies that differentially private learners need to pay an additional dimension d factor as well as a bit complexity of point

representation b factor over non-private learners. To the best of our knowledge such strong separation was not known before for problems defined over i.i.d. samples from a distribution (as opposed to worst case inputs). Note that this lower bound is also almost tight since $\log |\text{HS}_b^d| = O(d^2(\log d + b))$ [e.g. [Muroga, 1971](#)].

In the second result of [Theorem 1.2](#) we use the class Line_p of lines in \mathbb{Z}_p^2 (a plane over a finite field \mathbb{Z}_p). A lower bound on the one-way quantum communication complexity of this class was first given by [Aaronson \[2004\]](#) using his method based on a trace distance.

Finally, we consider PAC learning with (α, β) -differential privacy. Our lower bound of $\Omega(b)$ on SCDP of thresholds together with the upper bound of $O(16^{\log^*(b)} \cdot \log(1/\beta))$ from [[Beimel et al., 2013b](#)] immediately imply a separation between the sample complexities of pure and approximate differential privacy. We show a stronger separation for the concept class Line_p :

Theorem 1.3. *The sample complexity of (α, β) -differentially privately learning Line_p is $O(\frac{1}{\alpha} \log(1/\beta))$.*

Our upper bound is also simpler than the upper bound in [[Beimel et al., 2013b](#)]. See [Section 6](#) for details.

1.2 Related work

There is now an extensive amount of literature on differential privacy in machine learning and related areas which we cannot hope to cover here. The reader is referred to the excellent surveys in [[Sarwate and Chaudhuri, 2013](#), [Dwork and Roth, 2014](#)].

[Blum et al. \[2005\]](#) showed that algorithms that can be implemented in the statistical query (SQ) framework of [Kearns \[1998\]](#) can also be easily converted to differentially private algorithms. This result implies polynomial upper bounds on the sample (and computational) complexity of all learning problems that can be solved using statistical queries (which includes the vast majority of problems known to be solvable efficiently). Formal treatment of differentially private PAC and agnostic learning was initiated in the seminal work of [Kasiviswanathan et al. \[2011\]](#). Aside from the results we already mentioned, they separated SQ learning from differentially private learning. Further, they showed that SQ learning is (up to polynomial factors) equivalent to *local* differential privacy a more stringent model in which each data point is privatized before reaching the learning algorithm.

The results of this paper are for the distribution-independent learning, where the learner does not know the distribution over the domain. Another commonly-considered setting is distribution-specific learning in which the learner only needs to succeed with respect to a single fixed distribution \mathcal{D} known to the learner. Differentially private learning in this setting and its relaxation in which the learner only knows a distribution close to \mathcal{D} were studied by [Chaudhuri and Hsu \[2011\]](#). $\text{DRDim}_\varepsilon(C)$ restricted to a fixed distribution \mathcal{D} is denoted by $\text{DRDim}_\varepsilon^\mathcal{D}(C)$ and equals to the logarithm of the smallest ε -cover of C with respect to the disagreement metric given by \mathcal{D} (also referred to as the *metric entropy*). The standard duality between packing and covering numbers also implies that $\text{PRDim}_{\frac{\mathcal{D}}{2}, \delta}^\mathcal{D}(C) \geq \text{DRDim}_\varepsilon^\mathcal{D}(C) - \log(\frac{1}{1-\delta})$, and therefore these notions are essentially identical. It also follows from the prior work [[Kasiviswanathan et al., 2011](#), [Chaudhuri and Hsu, 2011](#)], that $\text{DRDim}_\varepsilon^\mathcal{D}(C)$ characterizes the complexity of differentially private PAC and agnostic learning up to the dependence on the error parameter ε in the same way as it does for (non-private) learning [[Benedek and Itai, 1991](#)]. Namely, $\Omega(\text{DRDim}_{2\varepsilon}^\mathcal{D}(C)/\alpha)$ samples are necessary to learn α -differentially privately with error ε (and even if only weaker label differential-privacy is desired [[Chaudhuri and Hsu, 2011](#)]) and $O(\text{DRDim}_{\varepsilon/2}^\mathcal{D}(C)/(\varepsilon\alpha))$ samples suffice for α -differentially private PAC learning. This implies that in this setting there are no dimension or bit-complexity costs incurred by differentially private learners. [Chaudhuri and Hsu \[2011\]](#) also show that doubling dimension at an appropriate scale can be used to give upper and lower bounds on sample complexity of distribution-specific private PAC learning that match up to logarithmic factors.

In a related problem of sanitization of queries from the concept class C the input is a database D of points in X and the goal is to output differentially privately a “synthetic” database \hat{D} such that for every

$f \in C$, $\left| \frac{1}{|D|} \sum_{x \in D} f(x) - \frac{1}{|\hat{D}|} \sum_{x \in \hat{D}} f(x) \right| \leq \varepsilon$. This problem was first considered by [Blum et al. \[2013\]](#) who showed an upper bound of $O(\text{VC}(C) \cdot \log(|X|))$ on the size of the database sufficient for this problem and also showed a lower bound of $\Omega(b)$ on the number of samples required for solving this problem when $X = I_b$ for $C = \text{Thr}_b$. It is easy to see that from the point of view of sample complexity this problem is at least as hard as (differentially private) *proper* agnostic learning of C [e.g. [Gupta et al., 2011](#)]. Therefore lower bounds on proper learning such as those in [[Beimel et al., 2010](#)] and [[Chaudhuri and Hsu, 2011](#)] apply to this problem and can be much larger than SCDP that we study. That said, to the best of our knowledge, the lower bound for linear threshold functions that we give was not known even for this harder problem. Aside from sample complexity this problem is also computationally intractable for many interesting classes C (see [[Ullman, 2013](#)] and references therein for recent progress).

Sample complexity of more general problems in statistics was investigated in several works starting with [Dwork and Lei \[2009\]](#) (measured alternatively via convergence rates of statistical estimators) [[Smith, 2011](#), [Chaudhuri and Hsu, 2012](#), [Duchi et al., 2013a,b](#)]. A recent work of [Duchi et al. \[2013a\]](#) shows a number of d -dimensional problems where differentially private algorithms must incur an additional factor d/α^2 cost in sample complexity. However their lower bounds apply only to a substantially more stringent local model of differential privacy and are known not to hold in the model we consider here.

Differentially private communication protocols were studied by [McGregor et al. \[2010\]](#) who showed that differential-privacy can be exploited to obtain a low-communication protocol and vice versa. Conceptually this result is similar to the characterization of sample complexity using PRDim given in [[Beimel et al., 2013a](#)]. Our contribution is orthogonal to [[McGregor et al., 2010](#)] since the main step in our work is going from a learning problem to a communication protocol for a different problem.

2 Preliminaries

2.1 Learning models

Definition 2.1. An algorithm A PAC learns a concept class C from n examples if for every $\epsilon > 0, \delta > 0, f \in C$ and distribution \mathcal{D} over X , A given access to $S = \{(x_i, \ell_i)\}_{i \in [n]}$ where each x_i is drawn randomly from \mathcal{D} and $\ell_i = f(x_i)$, outputs, with probability at least $1 - \delta$ over the choice of S and the randomness of A , a hypothesis h such that $\Pr_{x \sim \mathcal{D}}[f(x) \neq h(x)] \leq \epsilon$.

Agnostic learning: The *agnostic* learning model was introduced by [Haussler \[1992\]](#) and [Kearns et al. \[1994\]](#) in order to model situations in which the assumption that examples are labeled by some $f \in C$ does not hold. In its least restricted version the examples are generated from some unknown distribution P over $X \times \{0, 1\}$. The goal of an agnostic learning algorithm for a concept class C is to produce a hypothesis whose error on examples generated from P is close to the best possible by a concept from C . For a Boolean function h and a distribution P over $X \times \{0, 1\}$ let $\Delta(P, h) = \Pr_{(x, \ell) \sim P}[h(x) \neq \ell]$. Define $\Delta(P, C) = \inf_{h \in C} \{\Delta(P, h)\}$. [Kearns et al. \[1994\]](#) define agnostic learning as follows.

Definition 2.2. An algorithm A *agnostically* learns a concept class C if for every $\epsilon > 0, \delta > 0$, distribution P over $X \times \{0, 1\}$, A , given access to $S = \{(x_i, \ell_i)\}_{i \in [n]}$ where each (x_i, ℓ_i) is drawn randomly from P , outputs, with probability at least $1 - \delta$ over the choice of S and the randomness of A , a hypothesis h such that $\Delta(P, h) \leq \Delta(P, C) + \epsilon$.

In both PAC and agnostic learning model an algorithm that outputs a hypothesis in C is referred to as *proper*.

2.2 Differentially Private Learning

Two sample sets $S = \{(x_i, \ell_i)\}_{i \in [n]}$, $S' = \{(x'_i, \ell'_i)\}_{i \in [n]}$ are said to be *neighboring* if there exists $i \in [n]$ such that $(x_i, \ell_i) \neq (x'_i, \ell'_i)$, and for all $j \neq i$ it holds that $(x_j, \ell_j) = (x'_j, \ell'_j)$. For $\alpha, \beta > 0$, an algorithm A is (α, β) -differentially private if for all neighboring $S, S' \in (X \times \{0, 1\})^n$ and for all $T \subseteq \text{Range}(A)$:

$$\Pr[A(S) \in T] \leq e^\alpha \Pr[A(S') \in T] + \beta,$$

where the probability is over the randomness of A [Dwork et al., 2006]. When A is $(\alpha, 0)$ -differentially private we say that it satisfies *pure* differential privacy, which we also write as α -differential privacy.

Intuitively, each sample (x_i, ℓ_i) used by a learning algorithm is the record of one individual, and the privacy definition guarantees that by changing one record the output distribution of the learner does not change by much. We remark that, in contrast to the accuracy of learning requirement, the differential privacy requirement holds *in the worst case* for all neighboring sets of examples S, S' , not just those sampled i.i.d. from some distribution. We refer the reader to the literature for a further justification of this notion of privacy [Dwork et al., 2006].

The *sample complexity* $\text{SCDP}_{\alpha, \varepsilon, \delta}(C)$ is the minimal n such that it is information-theoretically possible to (ε, δ) -accurately and α -differentially privately PAC learn C with n examples. SCDP without subscripts refers to $\text{SCDP}_{1, \frac{1}{4}, \frac{1}{4}}$.

2.3 Representation Dimension

Definition 2.3 (Beimel et al., 2010). A class of functions H ε -represents C if for every $f \in C$ and every distribution \mathcal{D} over the input domain of f , there exists $h \in H$ such that $\Pr_{x \sim \mathcal{D}}[f(x) \neq h(x)] \leq \varepsilon$. The *deterministic representation dimension* of C , denoted as $\text{DRDim}_\varepsilon(C)$ equals $\log(|H|)$ for the smallest H that ε -represents C . We also let $\text{DRDim}(C) = \text{DRDim}_{\frac{1}{4}}(C)$.

Definition 2.4 (Beimel et al., 2013a). A distribution \mathcal{H} over sets of boolean functions on X is said to (ε, δ) -probabilistically represent C if for every $f \in C$ and distribution \mathcal{D} over X , with probability $1 - \delta$ over the choice of $H \stackrel{R}{\leftarrow} \mathcal{H}$, there exists $h \in H$ such that $\Pr_{x \sim \mathcal{D}}[h(x) \neq f(x)] \leq \varepsilon$. The (ε, δ) -probabilistic representation dimension $\text{PRDim}_{\varepsilon, \delta}(C)$ equals the minimal value of $\max_{H \in \text{supp}(\mathcal{H})} \log |H|$, where the minimum is over all \mathcal{H} that (ε, δ) -probabilistically represent C . We also let $\text{PRDim}(C) = \text{PRDim}_{\frac{1}{4}, \frac{1}{4}}(C)$.

Beimel et al. [2013a] proved the following characterization of SCDP by PRDim .

Theorem 2.5 (Kasiviswanathan et al., 2011, Beimel et al., 2013a).

$$\begin{aligned} \text{SCDP}_{\alpha, \varepsilon, \delta}(C) &= O\left(\frac{1}{\alpha \varepsilon} \left(\log(1/\varepsilon) \cdot \left(\text{PRDim}_{\frac{1}{4}, \frac{1}{4}}(C) + \log \log \frac{1}{\varepsilon \delta}\right) + \log \frac{1}{\delta}\right)\right). \\ \text{SCDP}_{\alpha, \varepsilon, \delta}(C) &= \Omega\left(\frac{1}{\alpha \varepsilon} \text{PRDim}_{1/4, 1/4}(C)\right). \end{aligned}$$

For agnostic learning we have that sample complexity is at most

$$O\left(\left(\frac{1}{\alpha \varepsilon} + \frac{1}{\varepsilon^2}\right) \left(\log(1/\varepsilon) \cdot \left(\text{PRDim}_{\frac{1}{4}, \frac{1}{4}}(C) + \log \log \frac{1}{\varepsilon \delta}\right) + \log \frac{1}{\delta}\right)\right).$$

This form of upper bound combines accuracy and confidence boosting from [Beimel et al., 2013a] to first obtain (ε, δ) -probabilistic representation and then the use of exponential mechanism as in [Kasiviswanathan et al., 2011]. The results in [Kasiviswanathan et al., 2011] show the extension of this bound to agnostic learning. Note that the characterization for PAC learning is tight up to logarithmic factors.

2.4 Communication Complexity

Let X and Y be some sets. A private-coin one-way protocol $\pi(x, y)$ from Alice who holds $x \in X$ to Bob who holds $y \in Y$ is given by Alice's randomized algorithm producing a communication σ and Bob's randomized algorithm which outputs a boolean value. We describe Alice's algorithm by a function $\pi_A(x; r_A)$ of the input x and random bits and Bob's algorithm $\pi_B(\sigma, y; r_B)$ by a function of input y , communication σ and random bits. (These algorithms need not be efficient.) The (randomized) output of the protocol on input (x, y) is the value of $\pi(x, y; r_A, r_B) \triangleq \pi_B(\pi_A(x; r_A), y; r_B)$ on a randomly and uniformly chosen r_A and r_B . The cost of the protocol $\text{CC}(\pi)$ is given by the maximum $|\sigma|$ over all $x \in X, y \in Y$ and all possible random coins.

A public-coin one-way protocol $\pi(x, y)$ is given by a randomized Alice's algorithm described by a function $\pi_A(x; r)$ and a randomized Bob's algorithm described by a function $\pi_B(\sigma, x; r)$. The (randomized) output of the protocol on input (x, y) is the value of $\pi(x, y; r) \triangleq \pi_B(\pi_A(x; r), y; r)$ on a randomly and uniformly chosen r . The cost of the protocol $\text{CC}(\pi)$ is defined as in the private-coin case.

Let $\Pi_\varepsilon^\rightarrow(g)$ denote the class of all private-coin one-way protocols π computing g with error ε , namely private-coin one-way protocols π satisfying for all $x \in X, y \in Y$

$$\Pr_{r_A, r_B} [\pi(x, y; r_A, r_B) = g(x, y)] \geq 1 - \varepsilon.$$

Define $\Pi_\varepsilon^{\rightarrow, \text{pub}}(g)$ similarly as the class of all public-coin one-way protocols π computing g . Define $R_\varepsilon^\rightarrow(g) = \min_{\pi \in \Pi_\varepsilon^\rightarrow(g)} \text{CC}(\pi)$ and $R_\varepsilon^{\rightarrow, \text{pub}}(g) = \min_{\pi \in \Pi_\varepsilon^{\rightarrow, \text{pub}}(g)} \text{CC}(\pi)$.

A deterministic one-way protocol π and its cost are defined as above but without dependence on random bits. We will also require distributional notions of complexity, where there is a fixed input distribution from which x, y are drawn. For a distribution μ over $X \times Y$, we define $\Pi_\varepsilon^\rightarrow(g; \mu)$ to be all deterministic one-way protocols π such that

$$\Pr_{(x, y) \sim \mu} [\pi(x, y) = g(x, y)] \geq 1 - \varepsilon.$$

Define $D_\varepsilon^\rightarrow(g; \mu) = \min_{\pi \in \Pi_\varepsilon^\rightarrow(g; \mu)} \text{CC}(\pi)$. A standard averaging argument shows that the quantity $D_\varepsilon^\rightarrow(g; \mu)$ remains unchanged even if we took the minimum over randomized (either public or private coin) protocols computing g with error $\leq \varepsilon$ (*i.e.* since there must exist a fixing of the private coins that achieves as good error as the average error).

Yao's minimax principle [Yao, 1977] states that for all functions g :

$$R_\varepsilon^{\rightarrow, \text{pub}}(g) = \max_{\mu} D_\varepsilon^\rightarrow(g; \mu). \quad (2.1)$$

Error in both public and private-coin protocols can be reduced by using several independent copies of the protocol and then taking a majority vote of the result. This implies that for every $\varepsilon, \gamma \in (0, 1/2)$,

$$R_\varepsilon^{\rightarrow, \text{pub}}(f) = O(R_{1/2-\gamma}^{\rightarrow, \text{pub}}(f) \cdot \log(1/\varepsilon)/\gamma^2). \quad (2.2)$$

Analogous statement holds for R^\rightarrow . This allows us to treat protocols with constant errors in $(0, 1/2)$ range as equivalent up to a constant factor in the communication complexity.

2.5 Littlestone's Dimension

While in this work we will not use the definition of the online mistake-bound model itself, we briefly describe it for completeness. In the online mistake-bound model learning proceeds in rounds. At the beginning of round t , a learning algorithm has some hypothesis h_t . In round t , the learner sees a point $x_t \in X$ and predicts $h_t(x_t)$. At the end of the round, the correct label y_t is revealed and the learner makes a mistake if $h_t(x_t) \neq y_t$. The learner then updates its hypothesis to h_{t+1} and this process continues. When learning a concept class C in

this model $y_t = f(x_t)$ for some unknown $f \in C$. The (sample) complexity of such learning is defined as the largest number of mistakes that any learning algorithm can be forced to make when learning C . Littlestone [1987] proved that it is exactly characterized by a dimension defined as follows.

Let C be a concept class over domain X . A *mistake tree* T over X and C is a binary tree in which each internal node v is labelled by a point $x_v \in X$ and each leaf ℓ is labelled by a concept $c_\ell \in C$. Further, for every node v and leaf ℓ : if ℓ is in the right subtree of v then $c_\ell(x_v) = 1$, otherwise $c_\ell(x_v) = 0$. We remark that a mistake tree over X and C does not necessarily include all concepts from C in its leaves. Such a tree is called *complete* if all its leaves are at the same depth. Littlestone's dimension $\text{LDim}(C)$ is defined as the depth of the deepest complete mistake tree T over X and C [Littlestone, 1987]. Littlestone's dimension is also known to exactly characterize the number of (general) equivalence queries required to learn C in Angluin's [1988] exact model of learning [Littlestone, 1987].

3 Equivalence between representation dimension and communication complexity

We relate communication complexity to private learning by considering the communication problem associated with evaluating a function f from a concept class C on an input $x \in X$. Formally, for a Boolean concept class C over domain X , define $\text{Eval}_C : C \times X \rightarrow \{0, 1\}$ to be the function defined as $\text{Eval}_C(f, x) = f(x)$. In a slight abuse of notation we use $\text{R}_{\varepsilon}^{\rightarrow, \text{pub}}(C)$ to denote $\text{R}_{\varepsilon}^{\rightarrow, \text{pub}}(\text{Eval}_C)$ (and similarly for $\text{R}_{\varepsilon}^{\rightarrow}(C)$).

Our main result is the following two bounds.

Theorem 3.1. *For any $\varepsilon \in [0, 1/2]$ and $\delta \in [0, 1]$, and any concept class C , it holds that:*

- $\text{PRDim}_{\varepsilon, \delta}(C) \leq \text{R}_{\varepsilon\delta}^{\rightarrow, \text{pub}}(C)$.
- $\text{PRDim}_{\varepsilon, \delta}(C) \geq \text{R}_{\varepsilon+\delta-\varepsilon\delta}^{\rightarrow, \text{pub}}(C)$.

Proof. (\leq): let π be the public-coin one-way protocol that achieves the optimal communication complexity c . For each choice of the public random coins r , let H_r denote the set of functions $h_\sigma(x) = \pi_B(\sigma, x; r)$ over all possible σ . Thus, each H_r has size at most 2^c . Let the distribution \mathcal{H} be to choose uniformly random r and then output H_r .

We show that this family (ε, δ) -probabilistically represents C . We know from the fact that π computes Eval_C with error $\varepsilon\delta$ that it must hold for all $f \in C$ and $x \in X$ that:

$$\Pr_r[\pi_B(\pi_A(f; r), x; r) \neq f(x)] \leq \varepsilon\delta.$$

In particular, it must hold for any distribution \mathcal{D} over X that:

$$\Pr_{x \sim \mathcal{D}, r}[\pi_B(\pi_A(f; r), x; r) \neq f(x)] \leq \varepsilon\delta.$$

Therefore, it must hold that

$$\Pr_r \left[\Pr_{x \sim \mathcal{D}}[\pi_B(\pi_A(f; r), x; r) \neq f(x)] > \varepsilon \right] < \delta.$$

Note that $\pi_B(\pi_A(f; r), x; r) \equiv h_{\pi_A(f; r)}(x) \in H_r$ and therefore, with probability $\geq 1 - \delta$ over the choice of $H_r \stackrel{R}{\leftarrow} \mathcal{H}$, there exists $h \in H_r$ such that $\Pr_{x \sim \mathcal{D}}[h(x) \neq f(x)] \leq \varepsilon$.

(\geq): let \mathcal{H} be the distribution over sets of boolean functions that achieves $\text{PRDim}_{\varepsilon,\delta}(C)$. We will show that for each distribution μ over inputs (f, x) , we can construct a $(\varepsilon + \delta - \varepsilon\delta)$ -correct protocol for Eval_C over μ that has communication bounded by $\text{PRDim}_{\varepsilon,\delta}(C)$. Namely, we will prove that

$$\max_{\mu} D_{\varepsilon+\delta-\varepsilon\delta}^{\rightarrow}(\text{Eval}_C; \mu) \leq \text{PRDim}_{\varepsilon,\delta}(C). \quad (3.1)$$

By Yao's minimax principle (Equation 2.1) [Yao, 1977] this implies that

$$R_{\varepsilon+\delta-\varepsilon\delta}^{\rightarrow, \text{pub}}(C) \leq \text{PRDim}_{\varepsilon,\delta}(C).$$

Fix μ . This induces a marginal distribution \mathcal{F} over functions $f \in C$ and for every $f \in C$ a distribution \mathcal{D}_f which is μ conditioned on the function being f (note that μ is equivalent to drawing f from \mathcal{F} and then x from \mathcal{D}_f). The protocol π is defined as follows: use public coins to sample $H \stackrel{R}{\leftarrow} \mathcal{H}$. Alice knows f and so knows the distribution \mathcal{D}_f . Alice sends the index of $h \in H$ such that $\Pr_{x \sim \mathcal{D}_f}[h(x) \neq f(x)] \leq \varepsilon$ if such h exists or an arbitrary $h \in H$ otherwise. Bob returns $h(x)$.

The error of this protocol can be analyzed as follows. Fix f and let G_f denote the event that $H \stackrel{R}{\leftarrow} \mathcal{H}$ contains h such that $\Pr_{x \sim \mathcal{D}_f}[h(x) \neq f(x)] \leq \varepsilon$. Observe that G_f is independent of \mathcal{D}_f so that even conditioned on G_f x remains distributed according to \mathcal{D}_f . Also, since $\mathcal{H}(\varepsilon, \delta)$ -probabilistically represents C , we know that for every f , $\Pr_r[G_f] \geq 1 - \delta$. Therefore we can then deduce that:

$$\begin{aligned} \Pr_{r,(f,x) \sim \mu} [\pi(f, x; r) = f(x)] &= \Pr_{r,(f,x) \sim \mu} [\pi(f, x; r) = f(x) \wedge G_f] + \Pr_{r,(f,x) \sim \mu} [\pi(f, x; r) = f(x) \wedge \neg G_f] \\ &\geq \Pr_{r,f \sim \mathcal{F}}[G_f] \cdot \Pr_{r,x \sim \mathcal{D}_f} [\pi(f, x; r) = f(x) \mid G_f] \\ &\geq (1 - \delta)(1 - \varepsilon) = 1 - \delta - \varepsilon + \varepsilon\delta. \end{aligned}$$

Thus π computes C with error at most $\varepsilon + \delta - \varepsilon\delta$ and it has communication bounded by $\text{PRDim}_{\varepsilon,\delta}(C)$. \blacksquare

We also establish an analogous equivalence for DRDim and private-coin protocols.

Theorem 3.2. *For any $\varepsilon \in [0, 1/2]$, it holds that:*

- $\text{DRDim}_{\varepsilon}(C) \leq R_{\varepsilon/2}^{\rightarrow}(C)$.
- $\text{DRDim}_{\varepsilon}(C) \geq R_{\varepsilon}^{\rightarrow}(C)$.

Proof. (\leq): let $R_{\varepsilon/2}^{\rightarrow}(C) = c$ and fix the private-coin one-way protocol π that achieves c . We define the deterministic representation H to be all functions $h_{\sigma}(x) = \text{maj}_{r_B}\{\pi(\sigma, x; r_B)\}$, i.e. the majority value of Bob's outputs on input x and communication σ . Observe that there are 2^c such functions (one for each σ possible) and therefore it suffices to show that H ε -deterministically represents C . To see this, observe that for each $f \in C$, and all $x \in X$, it holds that:

$$\Pr_{r_B, \sigma \stackrel{R}{\leftarrow} \pi_A(f; r_A)} [f(x) = \pi_B(\sigma, x; r_B)] \geq 1 - \varepsilon/2.$$

In particular, this means that for all distributions \mathcal{D} over X , it holds that

$$\Pr_{x \sim \mathcal{D}, r_B, \sigma \stackrel{R}{\leftarrow} \pi_A(f; r_A)} [f(x) = \pi_B(\sigma, x; r_B)] \geq 1 - \varepsilon/2.$$

By a standard averaging argument, there must exist at least one σ such that

$$\Pr_{x \sim \mathcal{D}, r_B} [f(x) = \pi_B(\sigma, x; r_B)] \geq 1 - \varepsilon/2.$$

Now say that x is bad if $\Pr_{r_B}[f(x) = \pi_B(\sigma, x; r_B)] < 1/2$. By the above, it follows that $\Pr_{x \sim \mathcal{D}}[x \text{ bad}] \leq \varepsilon$. By definition, if x is not bad then $f(x) = h_\sigma(x)$, since h_σ is the majority of $\pi_B(\sigma, x; r_B)$ over all r_B . Therefore

$$\Pr_{x \sim \mathcal{D}}[f(x) = h_\sigma(x)] \geq 1 - \varepsilon.$$

This implies that H ε -deterministically represents C .

(\geq): We first apply von-Neumann's Minimax theorem to the definition of deterministic representation. In particular, suppose H is the family of functions that achieves $\text{DRDim}_\varepsilon(C)$. Thus, for each $f \in C$ and each distribution \mathcal{D} over X , there exists $h \in H$ such that $\Pr_{\mathcal{D}}[h(x) = f(x)] \geq 1 - \varepsilon$. We define a zero-sum game for each f with the first player choosing a point $x \in X$ and the second player choosing a hypothesis $h \in H$ and the payoff of the second player being $|h(x) - f(x)|$. The definition of $\text{DRDim}_\varepsilon(C)$ implies that for every mixed strategy of the first player the second player has a pure strategy that achieves payoff of at least $1 - \varepsilon$. By the Minimax theorem there exists a distribution \mathfrak{h}_f over H such that, for every $x \in X$, it holds that

$$\mathbf{E}_{h \sim \mathfrak{h}_f}[|h(x) - f(x)|] = \Pr_{h \sim \mathfrak{h}_f}[h(x) = f(x)] \geq 1 - \varepsilon.$$

Our private-coin protocol π for Eval_C will be the following: on input f , Alice will use her private randomness to sample $h \sim \mathfrak{h}_f$ and send the index of h to Bob. Bob then outputs $h(x)$. Thus, for each f, x , it holds that

$$\Pr_{\pi}[\pi(f, x) = f(x)] = \Pr_{h \sim \mathfrak{h}_f}[h(x) = f(x)] \geq 1 - \varepsilon$$

and so the protocol computes Eval_C with error $\leq \varepsilon$. ■

An immediate corollary of these equivalences and eq.(2.2) is that $\text{DRDim}(C) = \Theta(\text{R}_{1/3}^{\rightarrow}(C))$ and $\text{PRDim}(C) = \Theta(\text{R}_{1/3}^{\rightarrow, \text{pub}}(C))$ as we stated in [Theorem 1.1](#).

3.1 Applications

Our equivalence theorems allow us to import many results from communication complexity into the context of private PAC learning, both proving new facts and simplifying proofs of previously known results in the process.

Separating SCDP and VC dimension. Define Thr_b as the family of functions $t_x : I_b \rightarrow \{0, 1\}$ for $x \in I_b$ where $t_x(y) = 1$ if and only if $y \geq x$. The lower bound follows from an observation that $\text{Eval}_{\text{Thr}_b}$ is equivalent to the “greater-than” function $\text{GT}_b(x, y) = 1$ if and only if $x > y$, where $x, y \in \{0, 1\}^b$ are viewed as binary representations of integers in I_b . Note $\text{Eval}_{\text{Thr}_b}(t_x, y) = 1 - \text{GT}_b(x, y)$ and therefore these functions are the same up to the negation. GT_b is a well studied function in communication complexity and it is known that $\text{R}_{1/3}^{\rightarrow, \text{pub}}(\text{GT}_b) = \Omega(b)$ [[Miltersen et al., 1998](#)]. By combining this lower bound with [Theorem 3.1](#) we obtain that $\text{VC}(\text{Thr}_b) = 1$ yet $\text{PRDim}(\text{Thr}_b) = \Omega(b)$. From [Theorem 2.5](#) it follows that $\text{SCDP}(\text{Thr}_b) = \Omega(b)$.

We note that it is known that VC dimension corresponds to the maximal distributional one-way communication complexity over all *product* input distributions. Hence this separation is analogous to separation of distributional one-way complexity over product distributions and the maximal distributional complexity over all distributions achieved using the greater-than function [[Kremer et al., 1999](#)].

We also give more such separations using lower bounds on PRDim based on Littlestone's dimension. These are discussed in [Section 4](#).

Accuracy and confidence boosting. Our equivalence theorems give a simple alternative way to reduce error in probabilistic and deterministic representations without using sequential boosting as was done in [Beimel et al., 2013a]. Given a private PAC learner with constant error, say $(\varepsilon, \delta) = (1/8, 1/8)$, one can first convert the learner to a communication protocol with error $1/4$, use $O(\log \frac{1}{\varepsilon'\delta'})$ simple independent repetitions (as in eq.(2.2)) to reduce the error to $\varepsilon'\delta'$, and then convert the protocol back into a (ε', δ') -probabilistic representation. The “magic” here happens when we convert between the communication complexity and probabilistic representation using min-max type arguments. This is the same tool that can be used to prove (computationally inefficient) boosting theorems.

Probabilistic vs. deterministic representation dimension. It was shown by Newman [1991] that public and private coin complexity are the same up to additive logarithmic terms. In our setting (and with a specific choice of error bounds to simplify presentation), Newman’s theorem implies that

$$R_{1/8}^{\rightarrow}(C) \leq R_{1/9}^{\rightarrow, \text{pub}}(C) + O(\log \log(|C||X|)). \quad (3.2)$$

We know by Sauer’s lemma that $\log |C| \leq O(\text{VC}(C) \cdot \log |X|)$, therefore we deduce that:

$$R_{1/8}^{\rightarrow}(C) \leq R_{1/9}^{\rightarrow, \text{pub}}(C) + O(\log \log \text{VC}(C) + \log \log |X|).$$

By our equivalence theorems, $\text{DRDim}(C) = \text{DRDim}_{1/4}(C) \leq R_{1/8}^{\rightarrow}(C)$ and $R_{1/9}^{\rightarrow, \text{pub}}(C) \leq \text{PRDim}_{1/16, 1/24}(C) = O(\text{PRDim}(C))$. This implies that

$$\text{DRDim}(C) = O(\text{PRDim}(C) + \log \log |X|).$$

A version of this was first proved in [Beimel et al., 2013a], whose proof is similar in spirit to the proof of Newman’s theorem. We also remark that the fact that $\text{DRDim}_{1/3}(\text{Point}_b) = \Omega(\log b)$ while $\text{PRDim}_{1/3}(\text{Point}_b) = O(1)$ [Beimel et al., 2010, 2013a] corresponds to the fact that the private-coin complexity of the equality function is $\Omega(\log b)$, while the public-coin complexity is $O(1)$. Here Point_b is the family of point functions over $\{0, 1\}^b$, *i.e.* functions that are zero everywhere except on a single point.

Simpler learning algorithms. Using our equivalence theorems, we can “import” results from communication complexity to give simple private PAC learners. For example, the well-known constant communication equality protocol using inner-product-based hashing can be converted to a probabilistic representation using Theorem 3.1, which can then be used to learn point functions. The resulting learning algorithm is somewhat simpler than the constant sample complexity learner for Point_b described in [Beimel et al., 2010] and we believe that this view also provides useful intuition. We remark that the probabilistic representation for Point_b that results from the communication protocol is known and was used for learning Point_b by Feldman [2009] in the context of evolvability. A closely related representation is also mentioned in [Beimel et al., 2013a].

Furthermore in some cases this connection can lead to *efficient* private agnostic learning algorithms. Namely, if there is a communication protocol for Eval_C where Bob’s algorithm is polynomial-time then one can run the exponential mechanism in time $2^{O(\text{PRDim}(C))}$ to differentially privately agnostically learn C .

4 Lower Bounds via Littlestone’s Dimension

In this section, we show that Littlestone’s dimension lower bounds the sample complexity of differentially private learning. Let C be a concept class over X of $\text{LDim } d$. Our proof is based on a reduction from the communication complexity of Eval_C to the communication complexity of Augmented Index problem on d bits. AugIndex is the promise problem where Alice gets a string $x_1, \dots, x_d \in \{0, 1\}^d$ and Bob gets $i \in [d]$

and x_1, \dots, x_{i-1} , and $\text{AugIndex}(x, (i, x_{[i-1]})) = x_i$ where $x_{[i-1]} = (x_1, \dots, x_{i-1})$. A variant of this problem in which the length of the prefix is not necessarily i but some additional parameter m was first explicitly defined by [Bar-Yossef et al. \[2004\]](#) who proved that it has randomized one-way communication complexity of $\Omega(d - m)$. The version defined above is from [\[Ba et al., 2010\]](#) where it is also shown that a lower bound for AugIndex follows from an earlier work of [\[Miltersen et al., 1998\]](#). We use the following lower bound for AugIndex .

Lemma 4.1. $R_{\varepsilon}^{\rightarrow}(\text{AugIndex}) \geq (1 - H(\varepsilon))d$, where $H(\varepsilon) = \varepsilon \log(1/\varepsilon) + (1 - \varepsilon) \log(1/(1 - \varepsilon))$ is the binary entropy function.

A proof of this lower bound can be easily derived by adapting the proof in [\[Bar-Yossef et al., 2004\]](#) and we include it in [Section A](#).

We now show that if $\text{LDim}(C) = d$ then one can reduce AugIndex on d bit inputs to Eval_C .

Lemma 4.2. *Let C be a concept class over X and $d = \text{LDim}(C)$. There exist two mappings $m_C : \{0, 1\}^d \rightarrow C$ and $m_X : \bigcup_{i \in [d]} \{0, 1\}^i \rightarrow X$ such that for every x and $i \in [d]$, the value of $m_C(x)$ on point $m_X(x_{[i-1]})$ is equal to $\text{AugIndex}(x, (i, x_{[i-1]})) = x_i$.*

Proof. By the definition of LDim , there exists a complete mistake tree T over X and C of depth d . Recall that a mistake tree over X and C is a binary tree in which each internal node is labelled by a point in X and each leaf is labelled by a concept in C . For $x \in \{0, 1\}^d$ consider a path from the root of the tree such that at step $j \in [d]$ we go to the left subtree if $x_j = 0$ and the right subtree if $x_j = 1$. Such path will end in a leaf which we denote by ℓ_x and the concept that labels it by c_x . For a prefix $x_{[i-1]}$, let $v_{x_{[i-1]}}$ denote the internal node at depth i on this path (with v_{\emptyset} being the root) and let $z_{x_{[i-1]}}$ denote the point in X which labels $v_{x_{[i-1]}}$.

We define the mapping m_C as $m_C(x) = c_x$ for all $x \in \{0, 1\}^d$ and the mapping m_X as $m_X(y) = z_y$ for all $y \in \bigcup_{i \in [d]} \{0, 1\}^i$. By the definition of a mistake tree over X and C , the value of the concept c_x on the point $z_{x_{[i-1]}}$ is determined by whether the leaf ℓ_x is in the right (1) or the left (0) subtree of the node $v_{x_{[i-1]}}$. Recall that the turns in the path from the root of the tree to ℓ_x are defined by the bits of x . At the node $v_{x_{[i-1]}}$, x_i determines whether ℓ_x will be in the right or the left subtree. Therefore $c_x(z_{x_{[i-1]}}) = x_i$. Therefore the mapping we defined reduces AugIndex to Eval_C . \blacksquare

An immediate corollary of [Lemma 4.2](#) and [Lemma 4.1](#) is the following lower bound.

Corollary 4.3. *Let C be a concept class over X and $d = \text{LDim}(C)$. $R_{\varepsilon}^{\rightarrow}(C) \geq (1 - H(\varepsilon))d$.*

A stronger form of this lower bound was proved by [Zhang \[2011\]](#) who showed that the power of Partition Tree lower bound technique for one-way *quantum* communication complexity of [Nayak \[1999\]](#) can be expressed in terms of LDim of the concept class associated with the communication problem.

4.1 Applications

We can now use numerous known lower bounds for Littlestone's dimension of C to obtain lower bounds on sample complexity of private PAC learning. Here we list several examples of known results where $\text{LDim}(C)$ is (asymptotically) larger than the VC dimension of C .

1. $\text{LDim}(\text{Thr}_b) = b$ [\[Littlestone, 1987\]](#). $\text{VC}(\text{Thr}_b) = 1$.
2. Let BOX_b^d denote the class of all axis-parallel rectangles over $[2^b]^d$, namely all concepts $r_{s,t}$ for $s, t \in [2^b]^d$ defined as $r_{s,t}(x) = 1$ if and only if for all $i \in [d]$, $s_i \leq x_i \leq t_i$. $\text{LDim}(\text{BOX}_b^d) \geq b \cdot d$ [\[Littlestone, 1987\]](#). $\text{VC}(\text{BOX}_b^d) = d + 1$.

3. Let HS_b^d denote class of all linear threshold functions over $[2^b]^d$. $\text{LDim}(\text{HS}_b^d) = b \cdot d(d-1)/2$. This lower bound is stated in [Maass and Turán, 1994a]. We are not aware of a published proof and therefore a proof based on counting arguments in [Muroga, 1971] appears in Section B for completeness. $\text{VC}(\text{HS}_b^d) = d + 1$.
4. Let BALL_b^d denote class of all balls over $[2^b]^d$, that is all functions obtained by restricting a Euclidean ball in \mathbb{R}^d to $[2^b]^d$. Then $\text{LDim}(\text{BALL}_b^d) = \Omega(b \cdot d^2)$ [Maass and Turán, 1994b]. $\text{VC}(\text{BALL}_b^d) = d + 1$.

5 Separation of PRDim from LDim

We next consider the question of whether PRDim is equal to LDim. It turns out that the communication complexity literature [Zhang, 2011] already contains the following counter-example separating PRDim and LDim. Define:

$$\text{Line}_p = \{f : \mathbb{Z}_p^2 \rightarrow \{0, 1\} : \exists a, b \in \mathbb{Z}_p^2 \text{ s.t. } f(x, y) = 1 \text{ iff } ax + b = y\}$$

It is easy to see that $\text{LDim}(\text{Line}_p) = 2$ (an online learning algorithm only needs two different counterexamples to the constant 0 function to recover the unknown line). It was also shown [Aaronson, 2004] that the quantum one-way communication complexity of $\text{Eval}_{\text{Line}_p}$ is $\Theta(\log p)$. This already implies a separation between LDim and PRDim using Theorem 3.1 and the fact that quantum one-way communication lower-bounds randomized public-coin communication.

We give a new information-theoretic and simpler proof of Aaronson's result for randomized public-coin communication. We start with a brief review of basic notions from information theory.

5.1 Information theory background

We will use the convention of letting bold-face \mathbf{a}, \mathbf{b} denote random variables and regular type a, b denote particular values that those random variables may take.

Recall the following definitions of entropy (all logarithms are base 2):

$$\begin{aligned} \text{(Shannon entropy)} \quad H(\mathbf{x}) &= \sum_x \Pr[\mathbf{x} = x] \log \frac{1}{\Pr[\mathbf{x}=x]} \\ \text{(Rényi entropy or collision entropy)} \quad H_2(\mathbf{x}) &= \log \frac{1}{\sum_x \Pr[\mathbf{x} = x]^2} \\ \text{(Min-entropy)} \quad H_\infty(\mathbf{x}) &= \min_x \log \frac{1}{\Pr[\mathbf{x}=x]} \end{aligned}$$

Recall that for all random variables \mathbf{x} over some universe X , it holds that $\log |X| \geq H(\mathbf{x}) \geq H_2(\mathbf{x}) \geq H_\infty(\mathbf{x})$.

The conditional Shannon entropy is defined as $H(\mathbf{x} | \mathbf{y}) = \mathbb{E}_{y \leftarrow \mathbf{y}} [H(\mathbf{x} | \mathbf{y} = y)]$.

The (Shannon) mutual information is defined as $\mathbf{I}(\mathbf{x}; \mathbf{y} | \mathbf{z}) = H(\mathbf{x} | \mathbf{z}) - H(\mathbf{x} | \mathbf{yz})$. The mutual information satisfies the *chain rule*:

$$\mathbf{I}(\mathbf{x}; \mathbf{yz} | \mathbf{m}) = \mathbf{I}(\mathbf{x}; \mathbf{y} | \mathbf{m}) + \mathbf{I}(\mathbf{x}; \mathbf{z} | \mathbf{ym}).$$

The Kullback-Leibler divergence (also called relative entropy) is defined as:

$$D(\mathbf{x} \parallel \mathbf{x}') = \sum_x \Pr[\mathbf{x} = x] \log \frac{\Pr[\mathbf{x} = x]}{\Pr[\mathbf{x}' = x]}.$$

For two jointly distributed random variables \mathbf{xy} , let $\langle \mathbf{x} \rangle \langle \mathbf{y} \rangle$ denote independent samples from the marginal distributions of \mathbf{x} and \mathbf{y} . If conditioned on some event E , we write $(\langle \mathbf{x} \rangle \langle \mathbf{y} \rangle | E) = \langle \mathbf{x} | E \rangle \langle \mathbf{y} | E \rangle$.

Recall the following characterization of mutual information in terms of Kullback-Leibler divergence:

$$\mathbf{I}(\mathbf{x}; \mathbf{y}) = D(\mathbf{x}\mathbf{y} \parallel \langle \mathbf{x} \rangle \langle \mathbf{y} \rangle). \quad (5.1)$$

Lemma 5.1. *Let \mathbf{x}, \mathbf{y} be jointly distributed random variables. Suppose the support of \mathbf{y} has size 2^s . Then for every $t > 0$ it holds that:*

$$\Pr_{y \stackrel{R}{\leftarrow} \mathbf{y}} [\mathbf{H}_\infty(\mathbf{x} \mid \mathbf{y} = y) < \mathbf{H}_\infty(\mathbf{x}) - s - t] < 2^{-t}.$$

Proof. Let $k = \mathbf{H}_\infty(\mathbf{x})$. Say that y is *bad* if $\mathbf{H}_\infty(\mathbf{x} \mid \mathbf{y} = y) < k - s - t$. By Bayes' rule and the definition of min-entropy, it holds that for all bad y there exists x such that

$$\frac{2^{-k}}{\Pr[\mathbf{y} = y]} \geq \frac{\Pr[\mathbf{x} = x]}{\Pr[\mathbf{y} = y]} \geq \frac{\Pr[\mathbf{x} = x \wedge \mathbf{y} = y]}{\Pr[\mathbf{y} = y]} = \Pr[\mathbf{x} = x \mid \mathbf{y} = y] > 2^{-k+s+t}.$$

Therefore $\Pr[\mathbf{y} = y] < 2^{-s-t}$. This implies that:

$$\Pr[\mathbf{y} \text{ is bad}] = \sum_{y \text{ is bad}} \Pr[\mathbf{y} = y] < 2^s \cdot 2^{-s-t} = 2^{-t}$$

which proves the lemma. ■

Let \mathbf{x}, \mathbf{y} be random variables over a common universe X . Recall the following two equivalent definitions of statistical distance.

$$\Delta(\mathbf{x}, \mathbf{y}) = \frac{1}{2} \sum_{x \in X} |\Pr[\mathbf{x} = x] - \Pr[\mathbf{y} = x]| = \max_{f: X \rightarrow \{0,1\}} |\Pr[f(\mathbf{x}) = 1] - \Pr[f(\mathbf{y}) = 1]|.$$

Recall Pinsker's inequality:

$$\mathbf{Lemma 5.2.} \quad \Delta(\mathbf{x}, \mathbf{y}) \leq \sqrt{D(\mathbf{x} \parallel \mathbf{y})/2}.$$

5.2 Lower Bound on Communication Complexity of Line_p

To obtain a lower bound on $\text{PRDim}(\text{Line}_p)$ we prove that $R_{1/5}^{\rightarrow, \text{pub}}(\text{Eval}_{\text{Line}_p}) \geq \log p - O(1)$.

Theorem 5.3. $R_{1/5}^{\rightarrow, \text{pub}}(\text{Line}_p) \geq \log p - 7$.

Proof. We set $\varepsilon = 1/5$ and let $\gamma = 2(\frac{1}{2} - \frac{1}{p} - 2/5)^2$. For $p \leq 128$ the claim obviously holds so we can assume that $p > 110$. This implies that $\gamma \geq 2/121$ and $\log(1/\gamma) \leq 6$.

By the min-max principle, it suffices to exhibit an input distribution μ for which

$$D_\varepsilon^{\rightarrow}(\text{Eval}_{\text{Line}_p}; \mu) \geq \log p - 1 - \log \frac{1}{\gamma} \geq \log p - 7.$$

Let us consider the distribution μ that first samples $b \stackrel{R}{\leftarrow} \{0, 1\}$ and then outputs a sample from μ_b defined as follows. We describe Alice's function f by a pair $(a, b) \in \mathbb{Z}_p^2$ that defines a line. The distribution μ_0 outputs uniform and independent pairs $(a, b), (x, y) \stackrel{R}{\leftarrow} \mathbb{Z}_p^2$ while μ_1 outputs uniform $(a, b) \stackrel{R}{\leftarrow} \mathbb{Z}_p^2$ and then uniform (x, y) satisfying $ax + y = b$.

It is clear that $\Pr_{(a,b,x,y) \sim \mu_0} [ax + y = b] = 1/p$ while $\Pr_{(a,b,x,y) \sim \mu_1} [ax + y = b] = 1$. Therefore, it must be for any protocol π that computes $\text{Eval}_{\text{Line}_p}$ with overall error ε over μ , the following must hold:

$$\frac{1}{2} \left(\Pr_{(a,b,x,y) \sim \mu_0} [\pi((a, b), (x, y)) = 1] - \frac{1}{p} \right) + \frac{1}{2} \left(1 - \Pr_{(a,b,x,y) \sim \mu_1} [\pi((a, b), (x, y)) = 1] \right) \leq \varepsilon.$$

And therefore

$$\Pr_{(a,b,x,y) \sim \mu_1} [\pi((a,b), (x,y)) = 1] - \Pr_{(a,b,x,y) \sim \mu_0} [\pi((a,b), (x,y)) = 1] \geq 1 - \frac{1}{p} - 2\varepsilon. \quad (5.2)$$

We will show that this is impossible for any one-way protocol π that communicates less than $c = \log p - 1 - \log \frac{1}{\gamma} \geq \log p - 7$ bits. Fix any such π , and let m be the message sent by Alice.

Let us say that m is *good* if $H_\infty(\mathbf{a}, \mathbf{b} \mid \mathbf{m} = m) \geq \log p + \log \frac{1}{\gamma} = 2 \log p - c - 1$. Here, all random variables are sampled according to μ_1 . Then by [Lemma 5.1](#) it holds that:

$$\Pr_{(a,b,x,y) \sim \mu_1, m=\pi(a,b)} [m \text{ is not good}] = \Pr[H_\infty(\mathbf{a}, \mathbf{b} \mid \mathbf{m} = m) < 2 \log p - c - 1] \leq \frac{1}{2}.$$

We next claim that:

Claim 5.4. *For any good m , $H_2(\mathbf{x}, \mathbf{y} \mid \mathbf{m} = m) \geq 2 \log p - \gamma$ (with respect to the distribution μ_1).*

We first prove the theorem using this claim. It follows that $H(\mathbf{x}, \mathbf{y} \mid \mathbf{m} = m) \geq H_2(\mathbf{x}, \mathbf{y} \mid \mathbf{m} = m) \geq 2 \log p - \gamma$. Observe that:

$$\mathbf{I}(\mathbf{x}, \mathbf{y}; \mathbf{m} \mid \mathbf{m} \text{ is good}) = H(\mathbf{x}, \mathbf{y} \mid \mathbf{m} \text{ is good}) - H(\mathbf{x}, \mathbf{y} \mid \mathbf{m}, \mathbf{m} \text{ is good}) \leq \gamma.$$

By the divergence characterization of mutual information, this implies that

$$D(\mathbf{x}, \mathbf{y}, \mathbf{m} \mid \mathbf{m} \text{ is good} \parallel \langle \mathbf{x}, \mathbf{y} \rangle \langle \mathbf{m} \mid \mathbf{m} \text{ is good} \rangle) \leq \gamma. \quad (5.3)$$

On the other hand, observe that the distribution $(\langle \mathbf{x}, \mathbf{y} \rangle \langle \mathbf{m} \mid \mathbf{m} \text{ is good} \rangle)$ is identical to the distribution $(\mathbf{x}', \mathbf{y}', \mathbf{m}' \mid \mathbf{m}' \text{ is good})$ sampled according to the distribution μ_0 . (The definition of \mathbf{m}' good is the same as for \mathbf{m} , since the marginal distribution of both μ_0 and μ_1 on the variables $\mathbf{a}, \mathbf{b}, \mathbf{m}$ is identical.)

Therefore we have by Pinsker's inequality and [Equation 5.3](#) that:

$$\Delta((\mathbf{x}, \mathbf{y}, \mathbf{m} \mid \mathbf{m} \text{ is good}), (\mathbf{x}', \mathbf{y}', \mathbf{m}' \mid \mathbf{m}' \text{ is good})) \leq \sqrt{\gamma/2} < \frac{1}{2} - \frac{1}{p} - 2\varepsilon.$$

Thus, overall we have that:

$$\begin{aligned} \Delta((\mathbf{x}, \mathbf{y}, \mathbf{m}), (\mathbf{x}', \mathbf{y}', \mathbf{m}')) &\leq \Pr[\mathbf{m} \text{ is not good}] + \Delta((\mathbf{x}, \mathbf{y}, \mathbf{m} \mid \mathbf{m} \text{ is good}), (\mathbf{x}', \mathbf{y}', \mathbf{m}' \mid \mathbf{m}' \text{ is good})) \\ &< 1 - \frac{1}{p} - 2\varepsilon. \end{aligned}$$

However, since the output of Bob depends only on x, y, m , it therefore follows that the probability that Bob outputs 1 under μ_1 is less than the probability that Bob outputs 1 under μ_0 plus $1 - \frac{1}{p} - 2\varepsilon$, and this contradicts [Equation 5.2](#), proving the theorem. \blacksquare

Proof of Claim 5.4. The joint distribution $\mathbf{a}, \mathbf{b}, \mathbf{x}, \mathbf{y}, \mathbf{m}$ drawn from μ_1 can be viewed in the following order: first sample \mathbf{m} from the marginal distribution, then sample \mathbf{a}, \mathbf{b} conditioned on \mathbf{m} , and then sample \mathbf{x}, \mathbf{y} a random point conditioned on $\mathbf{a}\mathbf{x} + \mathbf{b} = \mathbf{y}$.

Conditioned on $\mathbf{m} = m$, consider $\mathbf{x}_1, \mathbf{y}_1$ and $\mathbf{x}_2, \mathbf{y}_2$ sampled independently as just described. There are two ways a collision can occur: either $\mathbf{a}_1, \mathbf{b}_1$ and $\mathbf{a}_2, \mathbf{b}_2$ collide or they do not. In the first case $(\mathbf{x}_1, \mathbf{y}_1) = (\mathbf{x}_2, \mathbf{y}_2)$ occurs with probability $1/p$, in the second with probability at most $1/p^2$ since there is at most one point where the two lines intersect.

We formalize this as follows:

$$\begin{aligned}
\Pr[(\mathbf{x}_1, \mathbf{y}_1) = (\mathbf{x}_2, \mathbf{y}_2) \mid \mathbf{m} = m] &= \Pr[(\mathbf{x}_1, \mathbf{y}_1) = (\mathbf{x}_2, \mathbf{y}_2) \wedge (\mathbf{a}_1, \mathbf{b}_1) = (\mathbf{a}_2, \mathbf{b}_2) \mid \mathbf{m} = m] \\
&\quad + \Pr[(\mathbf{x}_1, \mathbf{y}_1) = (\mathbf{x}_2, \mathbf{y}_2) \wedge (\mathbf{a}_1, \mathbf{b}_1) \neq (\mathbf{a}_2, \mathbf{b}_2) \mid \mathbf{m} = m] \\
&\leq \frac{1}{p} \Pr[(\mathbf{a}_1, \mathbf{b}_1) = (\mathbf{a}_2, \mathbf{b}_2) \mid \mathbf{m} = m] + \frac{1}{p^2} \\
&\leq \frac{1 + \gamma}{p^2}.
\end{aligned}$$

In the above we used the fact that $\Pr[(\mathbf{a}_1, \mathbf{b}_1) = (\mathbf{a}_2, \mathbf{b}_2) \mid \mathbf{m} = m] \leq \frac{\gamma}{p}$ because m is good.

Finally, $H_2(\mathbf{x}, \mathbf{y} \mid \mathbf{m} = m) \geq \log \frac{p^2}{1+\gamma} > 2 \log p - \gamma$. ■

6 Separating pure and (α, β) -differential privacy

We prove that it is possible to learn Line_p with (α, β) -differential privacy and (ε, δ) accuracy using $O(\frac{1}{\varepsilon\alpha} \log \frac{1}{\beta} \log \frac{1}{\delta})$ samples. This gives further evidence that it is possible to obtain much better sample complexity with (α, β) -differential privacy than pure differential privacy. Our separation is somewhat stronger than that implied by our lower bound for Thr_b and the upper bound of $O(16^{\log^*(b)})$ in [Beimel et al., 2013b] since for Line_p we are able to *match* the non-private sample complexity (when the privacy and accuracy parameters are constant²), even though, as mentioned in the previous section, randomized one-way communication complexity and therefore the SCDP of Line_p is asymptotically $\Theta(\log p)$. We note that our learner is not proper since in addition to lines it may output point functions and the all zero function.

Theorem 6.1. *For any prime p , any $\varepsilon, \delta, \alpha, \beta \in (0, 1/2)$, one can (ε, δ) -accurately learn Line_p with (α, β) -differential privacy using $O(\frac{1}{\varepsilon\alpha} \log \frac{1}{\beta} \log \frac{1}{\delta})$ samples.*

We prove this theorem in two steps: first we construct a learner with poor dependence on δ and then amplify using the exponential mechanism to obtain a learner with good dependence on δ .

6.1 A learner with poor dependence on δ

Lemma 6.2. *For any prime p , any $\varepsilon, \delta, \alpha, \beta \in (0, 1/2)$, it suffices to take $O(\frac{1}{\varepsilon} 2^{6/\delta} \cdot \frac{1}{\alpha} \log \frac{1}{\beta\delta})$ samples in order to (ε, δ) -learn Line_p with (α, β) -differential privacy.*

Proof. At a high level, we run the basic (non-private) learner based on VC-dimension $O(\frac{1}{\alpha} \log \frac{1}{\beta})$ times. We use the fact that Line_p is *stable* in that after a constant number of samples, with high probability there is a *unique* hypothesis that classifies the samples correctly. (This is simply because any two distinct points on a line define the line.) Therefore, in each of the executions of the non-private learner, we are likely to recover the same hypothesis. We can then release this hypothesis (α, β) -privately using the ‘‘Propose-Test-Release’’ framework.

The main challenge in implementing this intuition is to eliminate corner cases, where with roughly probability $1/2$ the sample set may contain two distinct positively labeled points and with probability $1/2$ only a single positively labeled point, as this would lead to unstable outputs. We do this by randomizing the *number* of samples we take.

Let t be a number of samples, to be chosen later. Given t samples $(x_1, y_1), \dots, (x_t, y_t)$, our *basic learner* will do the following:

²Formally a bound for constant β is uninformative since weak $1/\beta$ dependence is achievable by naive subsampling. In our case the dependence on $1/\beta$ is logarithmic and we can ignore this issue.

1. See if there exist two distinct samples $(x_i, y_i) \neq (x_j, y_j)$ that are both classified positively. If so, output the unique line defined by these points.
2. Otherwise, see if there exists any sample (x_i, y_i) classified positively. Output the point function that outputs 1 on (x_i, y_i) and zero elsewhere.
3. Otherwise, output the constant 0 hypothesis.

Our *overall learner* uses the basic learner as follows: first sample an integer k uniformly from the interval $[\log(\ln(3/2)/\varepsilon), \log(\ln(3/2)/\varepsilon) + 6/\delta]$ and set $t = 2^k$. Set $\ell = \max\{\frac{12}{\alpha} \ln \frac{2}{\beta\delta} + 13, 72 \ln \frac{4}{\delta}\}$. Set $n = t\ell$.

1. Take n samples and cut them into ℓ subsamples of size t , and run the basic learner on each of these.
2. Let the returned hypotheses be h_1, \dots, h_ℓ . Define $\text{freq}(h_1, \dots, h_\ell) = \text{argmax}_h |\{h_i = h \mid i \in [\ell]\}|$, *i.e.* the most frequently occurring hypothesis, breaking ties using lexicographical order. We define $\bar{h} = \text{freq}(h_1, \dots, h_\ell)$. Compute c to be the smallest number of h_i that must be changed in order to change the most frequently occurring hypothesis, *i.e.*

$$c = \min \{c \mid \exists h'_1, \dots, h'_\ell, \text{freq}(h'_1, \dots, h'_\ell) \neq \bar{h}, c = |\{i \mid h_i \neq h'_i\}|\}.$$

3. If $c + \Lambda(1/\alpha) > \frac{1}{\alpha} \ln \frac{1}{2\beta} + 1$ then output \bar{h} , otherwise output the constant 0 hypothesis.

Here, $\Lambda(1/\alpha)$ denotes the Laplace distribution, whose density function at point x equals $\alpha e^{-\alpha|x|}$. It is easy to check that adding $\Lambda(1/\alpha)$ to a sum of Boolean values renders that sum α -differentially private [Dwork et al., 2006].

We analyze the overall learner. Observe that once t is fixed, the basic learner is deterministic.

Privacy: we prove that the overall learner is (α, β) -differentially private. Consider any two neighboring inputs $x, x' \in (\mathbb{Z}_p^2 \times \{0, 1\})^n$. There are two cases:

- The most frequent hypothesis \bar{h} returned by running the basic learner on the ℓ subsamples of x, x' is the same. In this case, there are two possible outputs of the mechanism, either \bar{h} or the 0 hypothesis. Due to the fact that we decide between them using a count with Laplace noise and the count has sensitivity 1, the probability assigned to either output changes by at most a multiplicative $e^{-\alpha}$ factor between x, x' .
- The most frequent hypotheses are different. In this case $c = 1$ for both x, x' . The probability of *not* outputting 0 in either case is given by

$$\Pr[\Lambda(1/\alpha) > \frac{1}{\alpha} \ln \frac{1}{2\beta}] = \beta.$$

Otherwise, in both cases they output 0.

Accuracy: we now show that the overall learner (ε, δ) -PAC learns. We claim that:

Claim 6.3. *Fix any hidden line f and any input distribution \mathcal{D} . With probability $1 - \delta/2$ over the choice of t , there is a unique hypothesis with error $\leq \varepsilon$ that the basic learner will output with probability at least $2/3$ when given t independent samples from \mathcal{D} .*

Let us first assume this claim is true. Then it is easy to show that the overall learner (ε, δ) -learns: suppose we are in the $1 - \delta/2$ probability case where there is a unique hypothesis with error $\leq \varepsilon$ output by the basic learner. Then, by Chernoff, since $\ell \geq 72 \ln \frac{4}{\delta}$ it holds that with probability $1 - \delta/4$ at least $7/12$ fraction of the basic learner outputs will be this unique hypothesis. This means that the number of samples that must be modified to change the most frequent hypothesis is $c \geq \frac{\ell}{12}$. Therefore since $\ell \geq \frac{12}{\alpha} \ln \frac{1}{\beta\delta} + 13$, in this case the probability that the overall learner does not output this unique hypothesis is bounded by:

$$\Pr[c + \Lambda(\frac{1}{\alpha}) \leq \frac{1}{\alpha} \ln \frac{1}{2\beta} + 1] \leq \Pr[\Lambda(\frac{1}{\alpha}) < -\frac{1}{\alpha} \ln \frac{2}{\delta}] = \frac{\delta}{4}.$$

Thus the overall probability of not returning an ε -good hypothesis is at most δ .

Proof of Claim 6.3 Fix a concept f defined by a line given by $(a, b) \in \mathbb{Z}_p^2$ and any input distribution \mathcal{D} over \mathbb{Z}_p^2 .

Define the following events $\text{None}_t, \text{One}_t, \text{Two}_t$ parameterized by an integer $t > 0$ and defined over the probability space of drawing $(x_1, y_1), \dots, (x_t, y_t)$ independently from \mathcal{D} :

- None_t is the event that all of the (x_i, y_i) are not on the line (a, b) .
- One_t is the event that there exists some (x_i, y_i) on the line (a, b) , and furthermore for every other (x_j, y_j) on the line (a, b) is in fact equal to (x_i, y_i) .
- Two_t is the event that there exists distinct $(x_i, y_i) \neq (x_j, y_j)$ that are both on the line (a, b) .

Next we will show that with probability $1 - \delta/2$ over the choice of t , one of these three events has probability at least $2/3$, and then we show that this suffices to imply the claim.

Let $r = \Pr_{(x,y) \sim \mathcal{D}}[f(x, y) = 1]$, let $q_{x,y} = \Pr_{(x',y') \sim \mathcal{D}}[(x', y') = (x, y)]$, and let $q = \max_{(x,y) \in f^{-1}(1)} q_{x,y}$. We can characterize the probabilities of $\text{None}_t, \text{One}_t, \text{Two}_t$ in terms of r, q, t as follows:

$$\begin{aligned} \Pr[\text{None}_t] &= (1 - r)^t, \\ \Pr[\text{One}_t] &= \sum_{(x,y) \in f^{-1}(1)} ((1 - r + q_{x,y})^t - (1 - r)^t), \\ \Pr[\text{Two}_t] &= 1 - \Pr[\text{None}_t] - \Pr[\text{One}_t]. \end{aligned}$$

The characterizations for $\text{None}_t, \text{Two}_t$ are obvious. The characterization of One_t is exactly the probability over all $(x, y) \in f^{-1}(1)$ that all samples are either labeled 0 or equal (x, y) , excluding the event that they are all labeled 0.

From the above and by considering the (x, y) maximizing $q_{x,y}$, we have the following bounds:

$$\Pr[\text{None}_t] \geq 1 - rt, \tag{6.1}$$

$$\Pr[\text{One}_t] \geq (1 - r + q)^t - (1 - r)^t \geq 1 - (r - q)t - e^{-rt}, \tag{6.2}$$

$$\Pr[\text{Two}_t] \geq (1 - e^{-rt/2})(1 - e^{-(r-q)t/2}). \tag{6.3}$$

The first two follow directly from the fact that for all $x \in \mathbb{R}$ it holds that $1 - x \leq e^x$ and also for all $x \in [0, 1]$ and $y \geq 1$ it holds that $(1 - x)^y \geq 1 - xy$. Equation 6.3 follows from the following argument. Two_t contains the sub-event where there is at least one positive example in the first $t/2$ samples and a different positive example in the second $t/2$ samples. The probability of this sub-event is lower-bounded by $(1 - (1 - r)^{t/2})(1 - (1 - r + q)^{t/2}) \geq (1 - e^{-rt/2})(1 - e^{-(r-q)t/2})$.

t is good with high probability. Let us say that t is good for None_t if $t \leq \frac{1}{3r}$. We say t is good for One_t if $t \in [\frac{\ln 6}{r}, \frac{1}{6(r-q)}]$. We say t is good for Two_t if $t \geq \frac{2 \ln 6}{r-q}$. (It is possible that some of these events may be empty, but this does not affect our argument.) Using [Equation 6.1](#), [Equation 6.2](#) and [Equation 6.3](#), it is clear that if t is good for some event, then the probability of that event is at least $2/3$.

Let us say t is good if it is good for any one of None_t , One_t , Two_t . t is good means the following when viewed on the logarithmic scale:

$$\log t \in [0, \log \frac{1}{r} - \log 3] \cup [\log \frac{1}{r} + \log \ln 6, \log \frac{1}{r-q} - \log 6] \cup [\log \frac{1}{r-q} + \log(2 \ln 6), \infty).$$

But this means that t is bad on the logarithmic scale is equivalent to:

$$\log t \in (\log \frac{1}{r} - \log 3, \log \frac{1}{r} + \log \ln 6) \cup (\log \frac{1}{r-q} - \log 6, \log \frac{1}{r-q} + \log(2 \ln 6)). \quad (6.4)$$

Thus, for any r , there are at most 3 integer values of $\log t$ that are bad. But recall that $t = 2^k$ where k is uniformly chosen from $\{\log(\ln(3/2)/\varepsilon), \dots, \log(\ln(3/2)/\varepsilon) + 6/\delta\}$. Therefore the probability that $k = \log t$ is one of the bad values defined in [Equation 6.4](#) is at most $\delta/2$.

When t is good, basic learner outputs unique accurate hypothesis. To conclude, we argue that when t is good then the basic learner will output a unique hypothesis with error $\leq \varepsilon$ with probability $\geq 2/3$. This is obvious when t is good for Two_t , since whenever the basic learner sees two points on the line, it recovers the exact line. It is also easy to see that when t is good for None_t , the basic learner outputs the 0 hypothesis with probability $2/3$, and this has error at most ε since

$$2/3 \leq \Pr[\text{None}_t] = (1-r)^t \leq e^{-rt} \Rightarrow r \leq \ln(3/2)/t \leq \varepsilon.$$

It remains to argue that the basic learner outputs a unique hypothesis with error at most ε when t is good for One_t . Observe that we have actually set the parameters so that when t is good for One_t , it holds that:

$$\Pr[\text{One}_t \wedge \text{unique positive point is } (x_{\max}, y_{\max})] \geq 2/3, \quad (6.5)$$

where $(x_{\max}, y_{\max}) = \arg\max_{(x,y) \in f^{-1}(1)} q_{x,y}$. Therefore, for such t , the basic learner will output the point function that is positive on exactly (x_{\max}, y_{\max}) with probability at least $2/3$.

To show that this point function has error at most ε , it suffices to prove that

$$\Pr[f(x, y) = 1 \wedge (x, y) \neq (x_{\max}, y_{\max})] = r - q \leq \varepsilon.$$

From [Equation 6.5](#), we deduce that:

$$2/3 \leq (1-r+q)^t - (1-r)^t \leq e^{-(r-q)t} \Rightarrow r-q \leq \ln(3/2)/t \leq \varepsilon.$$

This concludes the proof. ■

Improving dependence on δ : We now improve the exponential dependence on $1/\delta$ in [Lemma 6.2](#) to prove [Theorem 6.1](#). We will use the algorithm of [Lemma 6.2](#) with $\delta = 1/2$ and accuracy $\varepsilon/2$ repeated $k = O(\log(1/\delta))$ times independently in order to construct a set H of k hypotheses. We then draw a fresh sample S of $O(\log(1/\delta)/(\varepsilon\alpha))$ examples and select one of the hypotheses based on their error on S using the exponential mechanism of [[McSherry and Talwar, 2007](#)]. This mechanism chooses a hypothesis from H with probability proportional to $e^{-\alpha \cdot \text{err}_S(h)^2}$, where $\text{err}_S(h) = |\{(x, \ell) \in S \mid h(x) \neq \ell\}|$. Simple analysis [e.g. [Kasiviswanathan et al., 2011](#), [Beimel et al., 2013a](#)] then shows that the selection mechanism is α -differentially private and outputs a hypothesis that has error of at most ε on \mathcal{D} with probability at least $1 - \delta$. Note that each of the k copies of the low-confidence algorithm and the exponential mechanism are run on disjoint sample sets and therefore there is no privacy loss from such composition. Hence the resulting algorithm is also (α, β) -differentially private. We include formal details in [Section C](#).

7 Conclusions and Open Problems

Our work continues the investigation of the costs of privacy in standard classification models initiated by [Kasiviswanathan et al. \[2011\]](#). Our main result is a new connection to communication complexity that provides a rich set of results and techniques developed in the past 30 years in communication complexity to the study of differentially private learning. Most notably, we show that tools from information theory can be used to resolve several fundamental questions about SCDP. Our connection relies on the characterization of SCDP using natural notions of representation dimension introduced by [Beimel et al. \[2013a\]](#). We remark, however, that our lower bounds can also be proved more directly without relying on the results in [\[Beimel et al., 2013a\]](#). Implicit in our lower bounds is a lower bound on the mutual information between random examples and the hypothesis output by the private learning algorithm. On the other hand, it is known and easy to show that α -differential privacy gives an upper bound of $O(\alpha \cdot n)$ on the mutual information between n data points and the output of the algorithm (see for example [\[Dwork et al., 2015\]](#)).

While we focus on differential privacy our lower bounds have immediate implications in other settings where information about the examples needs to be stored or transmitted for the purposes of classification, such as distributed computation, streaming and low-memory computation.

Our work also demonstrates that PAC learning with approximate differential privacy can be substantially more sample efficient than learning with pure differential privacy. However our understanding of classification with approximate differential privacy still has some major gaps. Most glaringly, we do not know whether the sample complexity of (α, β) -differentially private PAC learning is different from the VC dimension (up to a $\text{poly}(1/\alpha, \log(1/\beta))$ factor). Also the separation from the pure differential privacy holds only for PAC learning and we do not know if it is also true for agnostic learning.

Acknowledgements

We are grateful to Kobbi Nissim for first drawing our attention to the intriguing problem of understanding the relationship between probabilistic representation dimension and VC dimension, and for valuable discussions regarding the sample complexity of privately learning threshold functions. We thank Nina Balcan and Avrim Blum who brought up the relationship of our bounds for intervals in [Section 3.1](#) to those based on Littlestone’s dimension. Their insightful comments and questions have led to our result in [Theorem 1.2](#). We also thank Sasha Rakhlin and Sasha Sherstov for useful suggestions and references.

D.X. was supported in part by the French ANR Blanc program under contract ANR-12-BS02-005 (RDAM project), by NSF grant CNS-1237235, a gift from Google, Inc., and a Simons Investigator grant to Salil Vadhan.

References

- S. Aaronson. Limitations of quantum advice and one-way communication. In *IEEE Conference on Computational Complexity*, pages 320–332, 2004.
- D. Angluin. Queries and concept learning. *Machine Learning*, 2:319–342, 1988.
- K. D. Ba, P. Indyk, E. Price, and D. P. Woodruff. Lower bounds for sparse recovery. In *SODA*, pages 1190–1197, 2010.
- Z. Bar-Yossef, T. S. Jayram, R. Krauthgamer, and R. Kumar. The sketching complexity of pattern matching. In *APPROX-RANDOM*, pages 261–272, 2004.

- A. Beimel, S. P. Kasiviswanathan, and K. Nissim. Bounds on the sample complexity for private learning and private data release. In *TCC*, pages 437–454, 2010.
- A. Beimel, K. Nissim, and U. Stemmer. Characterizing the sample complexity of private learners. In *ITCS*, pages 97–110, 2013a.
- A. Beimel, K. Nissim, and U. Stemmer. Private learning and sanitization: Pure vs. approximate differential privacy. In *APPROX-RANDOM*, pages 363–378, 2013b.
- G. M. Benedek and A. Itai. Learnability with respect to fixed distributions. *Theoretical Computer Science*, 86(2):377 – 389, 1991.
- A. Blum, C. Dwork, F. McSherry, and K. Nissim. Practical privacy: the SuLQ framework. In *PODS*, pages 128–138, 2005.
- A. Blum, K. Ligett, and A. Roth. A learning theory approach to noninteractive database privacy. *J. ACM*, 60(2):12, 2013.
- K. Chaudhuri and D. Hsu. Sample complexity bounds for differentially private learning. In *COLT*, pages 155–186, 2011.
- K. Chaudhuri and D. Hsu. Convergence rates for differentially private statistical estimation. In *ICML*, 2012.
- J. C. Duchi, M. I. Jordan, and M. J. Wainwright. Local privacy and statistical minimax rates. In *FOCS*, pages 429–438, 2013a.
- J. C. Duchi, M. J. Wainwright, and M. I. Jordan. Local privacy and minimax bounds: Sharp rates for probability estimation. In *NIPS*, pages 1529–1537, 2013b.
- C. Dwork and J. Lei. Differential privacy and robust statistics. In *STOC*, pages 371–380, 2009.
- C. Dwork and A. Roth. The algorithmic foundations of differential privacy. *Foundations and Trends in Theoretical Computer Science*, 9(3-4):211–407, 2014.
- C. Dwork and A. Smith. Differential privacy for statistics: What we know and what we want to learn. *Journal of Privacy and Confidentiality*, 1(2):135–154, 2009.
- C. Dwork, F. McSherry, K. Nissim, and A. Smith. Calibrating noise to sensitivity in private data analysis. In *TCC*, pages 265–284, 2006.
- C. Dwork, V. Feldman, M. Hardt, T. Pitassi, O. Reingold, and A. Roth. Generalization in adaptive data analysis and holdout reuse. *CoRR*, abs/1506, 2015.
- V. Feldman. Robustness of evolvability. In *COLT*, pages 277–292, 2009.
- A. Gupta, M. Hardt, A. Roth, and J. Ullman. Privately releasing conjunctions and the statistical query barrier. In *STOC*, pages 803–812, 2011.
- D. Haussler. Decision theoretic generalizations of the PAC model for neural net and other learning applications. *Information and Computation*, 100(1):78–150, 1992. ISSN 0890-5401.
- S. P. Kasiviswanathan, H. K. Lee, K. Nissim, S. Raskhodnikova, and A. Smith. What can we learn privately? *SIAM J. Comput.*, 40(3):793–826, June 2011.

- M. Kearns. Efficient noise-tolerant learning from statistical queries. *Journal of the ACM*, 45(6):983–1006, 1998.
- M. Kearns, R. Schapire, and L. Sellie. Toward efficient agnostic learning. *Machine Learning*, 17(2-3):115–141, 1994.
- I. Kremer, N. Nisan, and D. Ron. On randomized one-round communication complexity. *Computational Complexity*, 8(1):21–49, 1999.
- N. Littlestone. Learning quickly when irrelevant attributes abound: a new linear-threshold algorithm. *Machine Learning*, 2:285–318, 1987.
- W. Maass and Turán. *How fast can a threshold gate learn?*, pages 381–414. MIT Press, 1994a.
- W. Maass and G. Turán. Algorithms and lower bounds for on-line learning of geometrical concepts. *Machine Learning*, 14(1):251–269, 1994b.
- A. McGregor, I. Mironov, T. Pitassi, O. Reingold, K. Talwar, and S. P. Vadhan. The limits of two-party differential privacy. In *FOCS*, pages 81–90, 2010.
- F. McSherry and K. Talwar. Mechanism design via differential privacy. In *FOCS*, pages 94–103, 2007.
- P. B. Miltersen, N. Nisan, S. Safra, and A. Wigderson. On data structures and asymmetric communication complexity. *J. Comput. Syst. Sci.*, 57(1):37–49, 1998.
- S. Muroga. *Threshold logic and its applications*. Wiley-Interscience, New York, 1971.
- A. Nayak. Optimal lower bounds for quantum automata and random access codes. In *FOCS*, pages 369–377, 1999.
- I. Newman. Private vs. common random bits in communication complexity. *Inf. Process. Lett.*, 39(2):67–71, 1991.
- A. D. Sarwate and K. Chaudhuri. Signal processing and machine learning with differential privacy: Algorithms and challenges for continuous data. *IEEE Signal Process. Mag.*, 30(5):86–94, 2013.
- A. Smith. Privacy-preserving statistical estimation with optimal convergence rates. In *STOC*, pages 813–822, 2011.
- J. Ullman. Answering $n^{2+O(1)}$ counting queries with differential privacy is hard. In *STOC*, pages 361–370, 2013.
- L. G. Valiant. A theory of the learnable. *Communications of the ACM*, 27(11):1134–1142, 1984.
- A. Yao. Probabilistic computations: Toward a unified measure of complexity. In *FOCS*, pages 222–227, 1977.
- S. Zhang. On the power of lower bound methods for one-way quantum communication complexity. In *ICALP (1)*, pages 49–60, 2011.

A Proof of Lemma 4.1

As before, for $\varepsilon \in [0, 1]$, we let $H(\varepsilon)$ denote the binary entropy of ε , namely the entropy of the Bernoulli random variable that equals 1 with probability ε . We recall Fano's inequality (for the Boolean case):

Lemma A.1 (Fano's inequality). *Let \mathbf{x}, \mathbf{y} be Boolean random variables such that $\Pr[\mathbf{x} = \mathbf{y}] \geq 1 - \varepsilon$. Then it holds that $H(\mathbf{x} \mid \mathbf{y}) \leq H(\varepsilon)$.*

We can now prove Lemma 4.1.

Proof. By the Yao's [1977] min-max principle ,

$$R_{\varepsilon}^{\rightarrow}(\text{AugIndex}) \geq D_{\varepsilon}^{\rightarrow}(\text{AugIndex}; \mu),$$

where μ is the input distribution that samples uniform $\mathbf{x} \stackrel{R}{\leftarrow} \{0, 1\}^d$ and uniform $\mathbf{i} \stackrel{R}{\leftarrow} [d]$.

Consider any deterministic protocol π computing AugIndex with error at most ε over μ , and suppose that π uses $\delta \cdot d$ communication. Let $(\mathbf{x}, \mathbf{i}) \sim \mu$. Then $\mathbf{I}(\pi_A(\mathbf{x}); \mathbf{x}) \leq \delta \cdot d$. By the chain rule for mutual information it follows that

$$\begin{aligned} \delta \cdot d &\geq \mathbf{I}(\pi_A(\mathbf{x}); \mathbf{x}) \\ &= \sum_{i=1}^d \mathbf{I}(\mathbf{x}_i; \pi_A(\mathbf{x}) \mid \mathbf{x}_1, \dots, \mathbf{x}_{i-1}) \\ &= \sum_{i=1}^d (\mathbf{H}(\mathbf{x}_i \mid \mathbf{x}_1, \dots, \mathbf{x}_{i-1}) - \mathbf{H}(\mathbf{x}_i \mid \pi_A(\mathbf{x}), \mathbf{x}_1, \dots, \mathbf{x}_{i-1})) \\ &= \sum_{i=1}^d (1 - \mathbf{H}(\mathbf{x}_i \mid \pi_A(\mathbf{x}), \mathbf{x}_1, \dots, \mathbf{x}_{i-1})). \end{aligned}$$

We therefore deduce that (for \mathbf{i} uniform over $[d]$):

$$1 - \mathbf{H}(\mathbf{x}_{\mathbf{i}} \mid \pi_A(\mathbf{x}), \mathbf{x}_1, \dots, \mathbf{x}_{\mathbf{i}-1}, \mathbf{i}) \leq \delta. \quad (1.1)$$

By Fano's Inequality, we know that if the probability of guessing $\mathbf{x}_{\mathbf{i}}$ given $\pi_A(\mathbf{x}), \mathbf{x}_1, \dots, \mathbf{x}_{\mathbf{i}-1}, \mathbf{i}$ (which is exactly Bob's input) is at least $1 - \varepsilon$, then $\mathbf{H}(\mathbf{x}_{\mathbf{i}} \mid \pi_A(\mathbf{x}), \mathbf{x}_1, \dots, \mathbf{x}_{\mathbf{i}-1}, \mathbf{i}) \leq H(\varepsilon)$. From this and Equation 1.1, we deduce that $\delta \geq 1 - H(\varepsilon)$. \blacksquare

B Ldim Lower Bound for Halfspaces

Recall that HS_b^d denotes the concept class of all halfspaces over I_b^d , where $I_b = \{0, 1, \dots, 2^b - 1\}$. Our proof is based on the technique used in [Muroga, 1971] to prove a lower bound of $2^{d(d-1)}$ on the total number of distinct halfspaces over $\{0, 1\}^d$. As a first step we prove the following simple lemma (for $b = 1$ it can also be found in [Muroga, 1971]).

Lemma B.1. *For an integer $b \geq 1$ let f be a halfspace over I_b^d . There exists a vector $w \in \mathbb{Z}^d$ and an integer $\theta \in \mathbb{Z}$ such that:*

- (w, θ) represents f , that is $f(x) = 1$ if and only if $w \cdot x \geq \theta$;
- for every two distinct $x, x' \in I_b^d$, $w \cdot x \neq w \cdot x'$.

We refer to such a representation of f as collision-free.

Proof. Let (w', θ') be any integer weight representation of f (such representation always exists for a halfspace over integer points). We first create a margin around the decision boundary by setting $w'' = 2^{d+1}w'$ and $\theta = 2^{d+1}\theta' - 2^d$. Note that (w'', θ) also represents f and, in addition, for every x , $|w'' \cdot x - \theta| \geq 2^d$. We now define for every $i \in [d]$, $w_i = w''_i + 2^i$. It is easy to see that $|w \cdot x - w'' \cdot x| \leq 2^d - 1$ and therefore (w, θ) represents f . Further, d least significant bits in the binary representation of $w \cdot x$ are exactly equal to x and therefore the second condition is also satisfied. ■

Let (w, θ) be some fixed collision-free representation of a halfspace f . By ordering the elements of I_b^d according to the value of $w \cdot x$ we obtain a strict order over the $(2^b)^d$ elements of I_b^d . Further any threshold function on this order is of the form $w \cdot x \geq \theta'$ for some $\theta' \in \mathbb{Z}$. We exploit this observation to embed threshold functions into halfspaces. We can then use the well-known fact that LDim of threshold functions over an interval of size 2^b is b .

Theorem B.2. $\text{LDim}(\text{HS}_b^d) = (d(d-1)/2 + 1) \cdot b$.

Proof. We construct a mistake tree T_d over HS_b^d and I_b^d inductively over the dimension d . For $d = 1$, HS_b^d includes all threshold functions on I_b and therefore we define T_1 is the complete binary tree representing the binary search on this interval. Note that the depth of this tree is b .

Now for $d \geq 2$, let T_{d-1} be the complete mistake tree over HS_b^{d-1} and I_b^{d-1} of depth $((d-1)(d-2)/2 + 1) \cdot b$ given by our inductive construction. For every leaf ℓ let $f_\ell \in \text{HS}_b^{d-1}$ be the halfspace labeling the leaf. Let (w', θ) be a collision-free representation of f_ℓ (arbitrarily chosen but fixed for every possible halfspace). Let $Z_\ell = \{w' \cdot y \mid y \in I_b^{d-1}\}$. The collision-free property of (w', θ) implies that $|Z_\ell| = |I_b|^{d-1} = 2^{b(d-1)}$. Let $z_0 < z_1 \dots < z_{|Z_\ell|-1}$ denote the elements of Z_ℓ ordered by value and for every $j \leq 2^{b(d-1)}$, let y^j denote the point y such that $w' \cdot y = z_j$. For every $z \in Z_\ell$ let $f_{\ell,z}$ be the halfspace over I_b^d defined by (w, θ) where, $w_d = \theta - z$ and $w_i = w'_i$ for all $i \leq d-1$. Clearly, $f_{\ell,z}$ restricted to the $d-1$ dimensional subcube $I_b^{d-1} \times \{0\}$ (that is points x in I_b^d for which $x_d = 0$) is equivalent to f_ℓ . When restricted to the $d-1$ dimensional subcube $I_b^{d-1} \times \{1\}$, $f_{\ell,z}$ is equivalent to $w' \cdot y \geq z_j$. Therefore, up to renaming of the points $y^j \rightarrow j$ and functions $f_{\ell,z_j} \rightarrow t_j$, $F_\ell = \{f_{\ell,z}\}_{z \in Z_\ell}$ restricted to $I_b^{d-1} \times \{1\}$ is identical to the class of linear thresholds on interval $I_{b(d-1)} = \{0, 1, \dots, 2^{b(d-1)}\}$. This means that there exists a complete mistake tree T_ℓ for F_ℓ over $I_b^{d-1} \times \{1\}$ of depth $b(d-1)$.

Let T_d be the mistake tree obtained by (a) starting with T_{d-1} ; (b) replacing points in I_b^{d-1} that label nodes by the corresponding points in $I_b^{d-1} \times \{0\}$; (c) replacing each leaf ℓ of T_{d-1} with T_ℓ . We claim that this is a complete mistake tree for HS_b^d over I_b^d of depth $(d(d-1)/2 + 1) \cdot b$. The fact that this tree is a complete binary tree of depth $(d(d-1)/2 + 1) \cdot b$ follows immediately from our construction since $((d-1)(d-2)/2 + 1) \cdot b + (d-1)b = (d(d-1)/2 + 1) \cdot b$. Now let ℓ' be a leaf of T_d labeled by some halfspace f_{ℓ,z_j} . Let v' be a node in T_d labeled by point x such that ℓ' is in the subtree of v' . If v' is a node derived from node v in T_{d-1} then, by definition, ℓ is a leaf in the subtree of v in T_{d-1} and v is labeled by y such that $x = y0$. By our construction, $f_{\ell,z_j}(y0) = f_\ell(y)$ and therefore $f_{\ell,z_j}(x) = 1$ if and only if ℓ is in the right subtree of v which is equivalent to ℓ' being in the right subtree of v' .

If v' is a node in T_ℓ , then $x \in I_b^{d-1} \times \{1\}$. On points in $I_b^{d-1} \times \{1\}$ the function f_{ℓ,z_j} corresponds to the threshold function t_j on the interval $I_{b(d-1)}$ and consistency with T_d follows from the properties of the binary search tree T_ℓ for threshold functions. ■

C Improving Dependence on δ in Theorem 6.1

We now improve the exponential dependence on $1/\delta$ in Lemma 6.2 to prove Theorem 6.1. We first introduce the exponential mechanism of McSherry and Talwar [2007]. For simplicity we only describe its restriction to

the learning setting. Let H be a hypothesis class and define the “quality score function” $q(S, h) = |\{(x, y) \in S \mid h(x) = y\}|$. The exponential mechanism for q with privacy α is the following: given an input $S \in (X \times \{0, 1\})^n$, output $h \in H$ according to the distribution $\text{EM}(S)$ given by

$$\Pr[\text{EM}(S) = h] \propto e^{\alpha q(S, h)/2}.$$

We use the following theorem about the exponential mechanism. Let $q_{\max}(S) = \max_{h \in H} q(S, h)$.

Theorem C.1 (McSherry and Talwar, 2007). *The exponential mechanism is α -differentially private. Furthermore, for all $S \in (X \times \{0, 1\})^n$ and all $t > 0$, it holds that*

$$\Pr[q(S, \text{EM}(S)) < q_{\max}(S) - t] \leq |H|e^{-\alpha t/2}.$$

We now finish the proof of [Theorem 6.1](#).

Proof of Theorem 6.1. We will use the algorithm of [Lemma 6.2](#) with $\delta = \Theta(1)$ in order to construct a small set of hypotheses from which we’ll then select one using the exponential mechanism. More precisely:

1. Set $k = \log(2/\delta)/\log(4/3)$. Run the algorithm of [Lemma 6.2](#) k times independently with fresh samples and with $(\frac{\varepsilon}{4}, \frac{1}{4})$ -accuracy and (α, β) -differential privacy. Call the resulting hypotheses $H = \{h_1, \dots, h_k\}$.
2. Sample $m = \frac{16}{\varepsilon\alpha} \log \frac{4k}{\delta}$ additional samples, call this set S . Use the exponential mechanism to output $h \in H$.

The mechanism is (α, β) -differentially private because samples used to produce h_1, \dots, h_k are learned with (α, β) differential privacy (notice each sample can only affect one of the h_i , therefore the privacy loss does not add up when considering the set of all hypotheses). Also, the samples used to pick $h \in H$ are used via the exponential mechanism, which is also α -differentially private.

To analyze the accuracy, observe that since each of the h_i is produced using an $(\frac{\varepsilon}{4}, \frac{1}{4})$ -accurate learner, therefore by independence of the executions and our choice of k , it holds with probability $\geq 1 - (1 - \frac{1}{4})^k \geq 1 - \frac{\delta}{2}$ that H contains some h that has error at most $\frac{\varepsilon}{4}$.

Next, for any $h \in H$ with $\Pr_{x \sim \mathcal{D}}[f(x) \neq h(x)] > \varepsilon$, observe that by a standard multiplicative Chernoff bound with probability $1 - e^{-\varepsilon m/12}$ over the choice of S it holds that

$$\Pr_{x \leftarrow S} [f(x) \neq h(x)] > \frac{1}{2} \Pr_{x \sim \mathcal{D}} [f(x) \neq h(x)] > \varepsilon/2. \quad (3.1)$$

Similarly, for any $h \in H$ with $\Pr_{x \sim \mathcal{D}}[f(x) \neq h(x)] \leq \varepsilon/4$, with probability $1 - e^{-\varepsilon m/8}$ over the choice of S it holds that:

$$\Pr_{x \leftarrow S} [f(x) \neq h(x)] < \frac{3}{2} \Pr_{x \sim \mathcal{D}} [f(x) \neq h(x)] \leq 3\varepsilon/8. \quad (3.2)$$

Therefore by a union bound, it holds with probability $1 - ke^{-\varepsilon m/12} > 1 - \frac{\delta}{4}$ that for all $h \in H$ with error greater than ε , it holds that $\Pr_{x \leftarrow S} [f(x) \neq h(x)] > \varepsilon/2$, and for all $h \in H$ with error at most $\varepsilon/4$, it holds that $\Pr_{x \leftarrow S} [f(x) \neq h(x)] < 3\varepsilon/8$.

This implies that with probability $1 - 3\delta/4$ over the probability of computing H and sampling S , it suffices to output some $h \in H$ such that $q(S, h) \geq \max_{h' \in H} q(S, h') - \varepsilon/8$. This is because we are in the case where H contains a hypothesis with error $\leq \varepsilon/4$, and therefore by [Equation 3.2](#) it holds that $\max_{h' \in H} q(S, h') > |S|(1 - 3\varepsilon/8)$ and therefore any such h will have error $\leq \varepsilon/2$ over S . By [Equation 3.1](#), we deduce that any such h must have error $\leq \varepsilon$ over \mathcal{D} .

By [Theorem C.1](#), the probability that the exponential mechanism outputs such a h is at least $1 - ke^{-\alpha \varepsilon m/16} \geq 1 - \delta/4$. Therefore the overall probability of outputting an ε -good hypothesis is at least $1 - \delta$. ■