Statistical-Query Lower Bounds via Functional Gradients

Surbhi Goel*1, Aravind Gollakota^{†2}, and Adam Klivans^{‡2}

¹Microsoft Research NYC ²Department of Computer Science, University of Texas at Austin

October 21, 2020

Abstract

We give the first statistical-query lower bounds for agnostically learning any non-polynomial activation with respect to Gaussian marginals (e.g., ReLU, sigmoid, sign). For the specific problem of ReLU regression (equivalently, agnostically learning a ReLU), we show that any statistical-query algorithm with tolerance $n^{-(1/\epsilon)^b}$ must use at least $2^{n^c}\epsilon$ queries for some constant b,c>0, where n is the dimension and ϵ is the accuracy parameter. Our results rule out general (as opposed to correlational) SQ learning algorithms, which is unusual for real-valued learning problems. Our techniques involve a gradient boosting procedure for "amplifying" recent lower bounds due to Diakonikolas et al. (COLT 2020) and Goel et al. (ICML 2020) on the SQ dimension of functions computed by two-layer neural networks. The crucial new ingredient is the use of a nonstandard convex functional during the boosting procedure. This also yields a best-possible reduction between two commonly studied models of learning: agnostic learning and probabilistic concepts.

^{*}Supported by the JP Morgan AI Fellowship.

[†]Supported by NSF awards AF-1909204, AF-1717896, and a UT Austin Provost's Fellowship.

[‡]Supported by NSF awards AF-1909204, AF-1717896, and the NSF AI Institute for Foundations of Machine Learning (IFML). Work done while visiting the Institute for Advanced Study, Princeton, NJ.

1 Introduction

In this paper we continue a recent line of research exploring the computational complexity of fundamental primitives from the theory of deep learning [GKK19, YS19, DKKZ20, YS20, DGK⁺20, FCG20]. In particular, we consider the problem of fitting a single nonlinear activation to a joint distribution on $\mathbb{R}^n \times \mathbb{R}$. When the nonlinear activation is ReLU, this problem is referred to as ReLU regression or agnostically learning a ReLU. When the nonlinear activation is sign and the labels are Boolean, this problem is equivalent to the well-studied challenge of agnostically learning a halfspace [KKMS08].

We consider arguably the simplest possible setting—when the marginal distribution is Gaussian—and give the first statistical-query lower bounds for learning broad classes of nonlinear activations. The statistical-query model is a well-studied framework for analyzing the sample complexity of learning problems and captures most known learning algorithms. For common activations such as ReLU, sigmoid, and sign, we give complementary upper bounds, showing that our results cannot be significantly improved.

Let \mathcal{H} be a function class on \mathbb{R}^n , and let \mathcal{D} be a labeled distribution on $\mathbb{R}^n \times \mathbb{R}$ such that the marginal on \mathbb{R}^n is $D = \mathcal{N}(0, I_n)$. We say that a learner learns \mathcal{H} under \mathcal{D} with error ϵ if it outputs a function f such that

$$\underset{(x,y) \sim \mathcal{D}}{\mathbb{E}}[f(x)y] \ge \max_{h \in \mathcal{H}} \underset{(x,y) \sim \mathcal{D}}{\mathbb{E}}[h(x)y] - \epsilon.$$

One can show that this loss captures 0-1 error in the Boolean case, as well as squared loss in the ReLU case whenever the learner is required to output a nontrivial hypothesis (i.e., a hypothesis with norm bounded below by some constant c > 0). (See Appendices E and F for details.)

For ReLU regression, we obtain the following exponential lower bound:

Theorem 1.1. Let \mathcal{H}_{ReLU} be the class of ReLUs on \mathbb{R}^n with unit weight vectors. Suppose that there is an SQ learner capable of learning \mathcal{H}_{ReLU} under \mathcal{D} with error ϵ using $q(n, \epsilon, \tau)$ queries of tolerance τ . Then for any ϵ , there exists $\tau = n^{-(1/\epsilon)^b}$ such that $q(n, \epsilon, \tau) \geq 2^{n^c} \epsilon$ for some 0 < b, c < 1/2. That is, a learner must either use tolerance smaller than $n^{-(1/\epsilon)^b}$ or more than $2^{n^c} \epsilon$ queries.

Prior work due to Goel et al. [GKK19] gave a quasipolynomial SQ lower bound (with respect to correlational queries) for ReLU regression when the learner is required to output a ReLU as its hypothesis.

For the sigmoid activation we obtain the following lower bound:

Theorem 1.2. Consider the above setup with \mathcal{H}_{σ} , the class of unit-weight sigmoid units on \mathbb{R}^n . For any ϵ , there exists $\tau = n^{-\Theta(\log^2 1/\epsilon)}$ such that $q(n, \epsilon, \tau) \geq 2^{n^c} \epsilon$ for some 0 < c < 1/2.

We are not aware of any prior work on the hardness of agnostically learning a sigmoid with respect to Gaussian marginals.

For the case of halfspaces, a result of Kalai et al. [KKMS08] showed that any halfspace can be agnostically learned with respect to Gaussian marginals in time and sample complexity $n^{O(1/\epsilon^4)}$, which was later improved to $n^{O(1/\epsilon^2)}$ [DKN10]. The only known hardness result for this problem is due to Klivans and Kothari [KK14] who gave a quasipolynomial

lower bound based on the hardness of learning sparse parity with noise. Here we give the first exponential lower bound:

Theorem 1.3. Consider the above setup with \mathcal{H}_{hs} , the class of unit-weight halfspaces on \mathbb{R}^n . For any ϵ , there exists $\tau = n^{-\Theta(1/\epsilon)}$ such that $q(n, \epsilon, \tau) \geq 2^{n^c} \epsilon$ for some fixed constant 0 < c < 1/2.

Since it takes $\Theta(1/\tau^2)$ samples to simulate a query of tolerance τ , our constraint on τ here can be interpreted as saying that to avoid the exponential query lower bound, one needs sample complexity at least $\Theta(1/\tau^2) = n^{\Theta(1/\epsilon)}$, nearly matching the upper bound of [KKMS08, DKN10].

These results are formally stated and proved in Section 5. More generally, we show in Section 6 that our results give superpolynomial SQ lower bounds for agnostically learning any non-polynomial activation. (See Appendix A for some discussion of subtleties in interpreting these bounds.)

A notable property of our lower bounds is that they hold for *general* statistical queries. As noted by several authors [APVZ14, VW19], proving SQ lower bounds for real-valued learning problems often requires further restrictions on the types of queries the learner is allowed to make (e.g., correlational or Lipschitz queries).

Another consequence of our framework is the first SQ lower bound for agnostically learning monomials with respect to Gaussian marginals. In contrast, for the realizable (noiseless) setting, recent work due to Andoni et al. [ADHV19] gave an attribute-efficient SQ algorithm for learning monomials. They left open the problem of making their results noise-tolerant. We show in Section 7 that in the agnostic setting, no efficient SQ algorithm exists.

Theorem 1.4. Consider the above setup with \mathcal{H}_{mon} , the class of multilinear monomials of degree at most d on \mathbb{R}^n . For any $\epsilon \leq \exp(-\Theta(d))$ and $\tau \leq \epsilon^2$, $q(n, \epsilon, \tau) \geq n^{\Theta(d)} \tau^{5/2}$.

Our Approach Our approach deviates from the standard template for proving SQ lower bounds and may be of independent interest. In almost all prior work, SQ lower bounds are derived by constructing a sufficiently large family of nearly orthogonal functions with respect to the underlying marginal distribution. Instead, we will use a reduction-based approach:

- We show that an algorithm for agnostically learning a single nonlinear activation ϕ can be used as a subroutine for learning depth-two neural networks of the form $\psi(\sum_i \phi(w^i \cdot x))$ where ψ is any monotone, Lipschitz activation. This reduction involves an application of functional gradient descent via the Frank–Wolfe method with respect to a (nonstandard) convex surrogate loss.
- We apply recent work due to [DKKZ20] and [GGJ⁺20] that gives SQ lower bounds for learning depth-two neural networks of the above form in the probabilistic concept model. For technical reasons, our lower bound depends on the norms of these depth-two networks, and we explicitly calculate them for ReLU and sigmoid.
- We prove that the above reduction can be performed using only statistical queries. To do so, we make use of some subtle properties of the surrogate loss and the functional gradient method itself.

Our reduction implies the following new relationship between two well-studied models of learning: if concept class \mathcal{C} is efficiently agnostically learnable, then the class of monotone, Lipschitz functions of linear combinations of \mathcal{C} is learnable in the *probabilistic concept* model due to Kearns and Schapire [KS94]. We cannot hope to further strengthen the conclusion to *agnostic* learnability of monotone, Lipschitz functions of combinations of \mathcal{C} : the concept class of literals *is* agnostically learnable, but we show exponential SQ lower bounds for agnostically learning the class of majorities of literals, i.e., halfspaces (see also [KK14]).

Related Work Several recent papers have considered the computational complexity of learning simple neural networks [Bac17, GKKT17, YS20, FCG20, KK14, LSSS14, SVWX17, VW19, GKK19, GGJ⁺20, DKKZ20]. The above works either consider one-layer neural networks (as opposed to learning single neurons), or make use of discrete distributions (rather than Gaussian marginals), or hold for narrower classes of algorithms (rather than SQ algorithms). Goel et al. [GKK19] give a quasipolynomial correlational SQ lower bound for proper agnostic learning of ReLUs with respect to Gaussian marginals. They additionally give a similar computational lower bound assuming the hardness of learning sparse parity with noise.

The idea of using functional gradient descent to learn one hidden layer neural networks appears in work due to Bach [Bac17], who considered an "incremental conditional gradient algorithm" that at each iteration implicitly requires an agnostic learner to complete a "Frank–Wolfe step." A key idea in our work is to optimize with respect to a particular convex functional (surrogate loss) in order to obtain SQ learnability for depth-two neural networks with a nonlinear output activation. We can then leverage SQ lower bounds for this broader class of neural networks.

Functional gradient descent or gradient boosting methods have been used frequently in learning theory, especially in online learning (see e.g., [Fri01, MBBF00, SF12, BHKL15, Haz16].)

For Boolean functions, the idea to use boosting to learn majorities of a base class appeared in Jackson [Jac97], who boosted a weak parity learning algorithm in order to learn thresholds of parities (TOP). Agnostic, distribution-specific boosting algorithms for Boolean functions have appeared in works due to Kalai and Kanade [KK09] and also Feldman [Fel10]. Agnostic boosting in the context of the SQ model is explored in [Fel12], where an SQ lower bound is given for agnostically learning monotone conjunctions with respect to the uniform distribution on the Boolean hypercube.

The SQ lower bounds we obtain for agnostically learning halfspaces can be derived using one of the above boosting algorithms due to Kalai and Kanade [KK09] or Feldman [Fel10] in place of functional gradient descent, as halfspaces are Boolean functions.

Independent Work Independently and concurrently, Diakonikolas et al. [DKZ20] have obtained similar results for agnostically learning halfspaces and ReLUs. Rather than using a reduction-based approach, they construct a hard family of Boolean functions. They show that an agnostic learner for halfspaces or ReLUs would yield a learner for this family, which would solve a hard unsupervised distribution-learning problem considered in [DKS17]. Quantitatively, the lower bound they obtain is that agnostic learning of halfspaces or ReLUs up to excess error ϵ using queries of tolerance $n^{-\text{poly}(1/\epsilon)}$ requires at least $n^{\text{poly}(1/\epsilon)}$ queries.

These results are technically incomparable with ours. For queries of similar tolerance, our bound of $2^{n^c}\epsilon$ scales exponentially with n whereas theirs only scales polynomially, so that for any constant ϵ our bound is exponentially stronger. But our bound does not scale directly with $1/\epsilon$ (other than via the induced constraint on tolerance, which does scale as $n^{-\text{poly}(1/\epsilon)}$). Our work also extends to general non-polynomial activations, while theirs does not.

Organization We cover the essential definitions, models and existing lower bounds that we need in the preliminaries. Our main reduction, which says that if we could agnostically learn a single neuron, then we could learn depth-two neural networks composed of such neurons, is set up as follows. In Section 3 we explain our usage of functional gradient descent, with Assumption 3.1 formally stating the kind of agnostic learning guarantee we require for a single neuron. The main reduction itself is Theorem 4.1, the subject of Section 4. In Sections 5, 6 and 7 we derive the formal lower bounds which follow as a consequence of our reduction. Finally in Section 8, we contrast these lower bounds by also including some simple upper bounds.

2 Preliminaries

Notation Let D be a distribution over \mathbb{R}^n , which for us will be the standard Gaussian $\mathcal{N}(0,I_n)$ throughout. We will work with the L^2 space $L^2(\mathbb{R}^n,D)$ of functions from \mathbb{R}^n to \mathbb{R} , with the inner product given by $\langle f,g\rangle_D=\mathbb{E}_D[fg]$. The corresponding norm is $\|f\|_D=\sqrt{\mathbb{E}_D[f^2]}$. We refer to the ball of radius R as $\mathcal{B}_D(R)=\{f\in L^2(\mathbb{R}^n,D)\mid \|f\|_D\leq R\}$. We will omit the subscripts when the meaning is clear from context. Given vectors $u,v\in\mathbb{R}^n$, we will refer to their Euclidean dot product by $u\cdot v$ and the Euclidean norm by $\|u\|_2$. Given a function $\ell(a,b)$ we denote its partial derivative with respect to its first parameter, $\frac{\partial \ell}{\partial a}(a,b)$, by $\partial_1 \ell(a,b)$.

A Boolean probabilistic concept, or p-concept, is a function that maps each point x to a random $\{\pm 1\}$ -valued label y in such a way that $\mathbb{E}[y|x] = f^*(x)$ for a fixed function $f^*: \mathbb{R}^n \to [-1,1]$, known as its conditional mean function. We will use D_{f^*} to refer to the (unique) induced labeled distribution on $\mathbb{R}^n \times \{\pm 1\}$, i.e. we say $(x,y) \sim D_{f^*}$ if the marginal distribution of x is D and $\mathbb{E}[y|x] = f^*(x)$. We also sometimes use $y \sim f^*(x)$ to say that $y \in \{\pm 1\}$ and $\mathbb{E}[y|x] = f^*(x)$.

Statistical Query (SQ) Model A statistical query is specified by a query function $\phi: \mathbb{R}^n \times \mathbb{R} \to [-1, 1]$. Given a labeled distribution \mathcal{D} on $\mathbb{R}^n \times \mathbb{R}$, the SQ model allows access to an SQ oracle (known as the STAT oracle in the SQ literature) that accepts a query ϕ of specified tolerance τ , and responds with a value in $[\mathbb{E}_{(x,y)\sim\mathcal{D}}[\phi(x,y)]-\tau,\mathbb{E}_{(x,y)\sim\mathcal{D}}[\phi(x,y)]+\tau]$. One can interpret the tolerance τ as capturing the notion of sample complexity in traditional PAC algorithms. Specifically, it takes $\Theta(1/\tau^2)$ samples to simulate a query of tolerance τ , and this is sometimes referred to as the estimation complexity of an SQ algorithm.

Let \mathcal{C} be a class of Boolean p-concepts over \mathbb{R}^n , and let D be a distribution on \mathbb{R}^n . We say that a learner learns \mathcal{C} with respect to D up to L^2 error ϵ if, given only SQ oracle access to D_{f^*} for some unknown $f^* \in \mathcal{C}$, and using arbitrary queries, it is able to output $f: \mathbb{R}^n \to [-1, 1]$ such that $||f - f^*||_D \le \epsilon$. It is worth emphasizing that a query to D_{f^*} takes in a Boolean rather than a real-valued label, i.e. is really of the form $\phi : \mathbb{R}^n \times \{\pm 1\} \to [-1, 1]$. In contrast, a query to a generic distribution \mathcal{D} on $\mathbb{R}^n \times \mathbb{R}$ takes in real-valued labels, and in Assumption 3.1 we define a form of learning that operates in this more generic setting.

One of the chief features of the SQ model is that one can give strong information theoretic lower bounds on learning a class \mathcal{C} in terms of its so-called statistical dimension.

Definition 2.1. Let D be a distribution on \mathbb{R}^n , and let \mathcal{C} be a real-valued or Boolean concept class on \mathbb{R}^n . The average (un-normalized) correlation of \mathcal{C} is defined to be $\rho_D(\mathcal{C}) = \frac{1}{|\mathcal{C}|^2} \sum_{c,c' \in \mathcal{C}} |\langle c,c' \rangle_D|$. The statistical dimension on average at threshold γ , $\mathrm{SDA}_D(\mathcal{C},\gamma)$, is the largest d such that for all $\mathcal{C}' \subseteq \mathcal{C}$ with $|\mathcal{C}'| \geq |\mathcal{C}|/d$, $\rho_D(\mathcal{C}') \leq \gamma$.

In the p-concept setting, lower bounds against general queries in terms of SDA were first formally shown in $[GGJ^+20]$.

Theorem 2.2 ([GGJ⁺20], Cor. 4.6). Let D be a distribution on \mathbb{R}^n , and let \mathcal{C} be a p-concept class on \mathbb{R}^n . Say our queries are of tolerance τ , the final desired L^2 error is ϵ , and that the functions in \mathcal{C} satisfy $||f^*|| \geq \beta$ for all $f^* \in \mathcal{C}$. For technical reasons, we will require $\tau \leq \epsilon^2$, $\epsilon \leq \beta/3$ (see Appendix A for some discussion). Then learning \mathcal{C} up to L^2 error ϵ (we may pick ϵ as large as $\beta/3$) requires at least $\mathrm{SDA}_D(\mathcal{C}, \tau^2)$ queries of tolerance τ .

A recent result of Diakonikolas et al [DKKZ20] gave the following construction of one-layer neural networks on \mathbb{R}^n with k hidden units, i.e. functions of the form $g(x) = \psi(\sum_{i=1}^k a_i \phi(x \cdot w_i))$ for activation functions $\psi, \phi : \mathbb{R} \to \mathbb{R}$ and weights $w_i \in \mathbb{R}^n, a_i \in \mathbb{R}$.

Theorem 2.3 ([DKKZ20]). There exists a class \mathcal{G} of one-layer neural networks on \mathbb{R}^n with k hidden units such that for some universal constant 0 < c < 1/2 and $\gamma = n^{\Theta(k(c-1/2))}$, $SDA(\mathcal{G}, \gamma) \geq 2^{n^c}$. This holds for any $\psi : \mathbb{R} \to [-1, 1]$ that is odd, and $\phi \in L^2(\mathbb{R}, \mathcal{N}(0, 1))$ that has a nonzero Hermite coefficient of degree greater than k/2. Further, the weights satisfy $|a_i| = 1/k$ and $||w_i||_2 = 1$ for all i.

We will be interested in the following special cases. Full details of the construction and proofs of the norm lower bounds are in Appendix B.

Corollary 2.4. For the following instantiations of \mathcal{G} , with accompanying norm lower bound β (i.e. such that $||g|| \geq \beta$ for all $g \in \mathcal{G}$), there exist $\tau = n^{-\Theta(k)}$ and $\epsilon \geq \tau$ such that learning \mathcal{G} up to L^2 error ϵ requires at least 2^{n^c} queries of tolerance τ , for some 0 < c < 1/2.

- (a) ReLU nets: $\psi = \tanh$, $\phi = \text{ReLU}$. Then $\beta = \Omega(1/k^6)$ (Lemma B.4), so we may take $\epsilon = \Theta(1/k^6)$.
- (b) Sigmoid nets: $\psi = \tanh$, $\phi = \sigma$. Then $\beta = \exp(-O(\sqrt{k}))$ (Lemma B.6), so we may $take \ \epsilon = \exp(-\Theta(\sqrt{k}))$.
- (c) Majority of halfspaces: $\psi = \phi = \text{sign}$. Being Boolean functions, here $\beta = 1$ exactly, so we may take $\epsilon = \Theta(1)$.

Convex Optimization Basics Over a general inner product space \mathcal{Z} , a function $p: \mathcal{Z} \to \mathbb{R}$ is convex if for all $\alpha \in [0,1]$ and $z,z' \in \mathcal{Z}$, $p(\alpha z + (1-\alpha)z') \leq \alpha p(z) + (1-\alpha)p(z')$.

We say that $s \in \mathcal{Z}$ is a subgradient of p at z if $p(z+h)-p(z) \geq \langle s,h \rangle$. We say that p is β -smoothly convex if for all $z,h \in \mathcal{Z}$ and any subgradient s of p at z,

$$p(z+h) - p(z) - \langle s, h \rangle \le \frac{\beta}{2} ||h||^2.$$

If there is a unique subgradient of p at z, we simply refer to it as the gradient $\nabla p(z)$. It is easily proven that smoothly convex functions have unique subgradients at all points. Another standard property is the following: for any $z, z' \in \mathcal{Z}$,

$$p(z) - p(z') \le \langle \nabla p(z), z - z' \rangle - \frac{1}{2\beta} \|\nabla p(z) - \nabla p(z')\|^2. \tag{1}$$

In this paper we will be concerned with convex optimization using the Frank–Wolfe variant of gradient descent, also known as conditional gradient descent. In order to eventually apply this framework to improper learning, we will consider a slight generalization of the standard setup. Let $\mathcal{Z}' \subset \mathcal{Z}$ both be compact, convex subsets of our generic inner product space. Say we have a β -smoothly convex function $p: \mathcal{Z} \to \mathbb{R}$, and we want to solve $\min_{z \in \mathcal{Z}'} p(z)$, i.e. optimize over the smaller domain, while allowing ourselves the freedom of finding subgradients that lie in the larger \mathcal{Z} . The Frank–Wolfe algorithm in this "improper" setting is Algorithm 1.

Algorithm 1 Frank-Wolfe gradient descent over a generic inner product space

```
Start with an arbitrary z_0 \in \mathcal{Z}.

for t = 0, ..., T do

Let \gamma_t = \frac{2}{t+2}.

Find s \in \mathcal{Z} such that \langle s, -\nabla p(z_t) \rangle \geq \max_{s' \in \mathcal{Z}'} \langle s', -\nabla p(z_t) \rangle - \frac{1}{2} \delta \gamma_t C_p.

Let z_{t+1} = (1 - \gamma_t) z_t + \gamma_t s.

end for
```

The following theorem holds by standard analysis (see e.g. [Jag13]). For convenience, we provide a self-contained proof in Appendix D.

Theorem 2.5. Let $\mathcal{Z}' \subseteq \mathcal{Z}$ be convex sets, and let $p : \mathcal{Z} \to \mathbb{R}$ be a β -smoothly convex function. Let $C_p = \beta \operatorname{diam}(\mathcal{Z})^2$. For every t, the iterates of Algorithm 1 satisfy

$$p(z_t) - \min_{z' \in \mathcal{Z}'} p(z') \le \frac{2C_p}{t+2} (1+\delta).$$

3 Functional gradient descent

Let $\ell: \mathbb{R} \times \mathbb{R} \to \mathbb{R}$ be a loss function. Given a p-concept f^* and its corresponding labeled distribution D_{f^*} , the population loss of a function $f: \mathbb{R}^n \to \mathbb{R}$ is given by $L(f) = \mathbb{E}_{(x,y) \sim D_{f^*}}[\ell(f(x),y)]$. We will view L as a mapping from $L^2(\mathbb{R}^n,D)$ to \mathbb{R} , and refer to it as the loss functional. The general idea of functional gradient descent is to try to find an f in a class of functions \mathcal{F} that minimizes L(f) by performing gradient descent in function space. When using Frank-Wolfe gradient descent, the key step in every iteration is to find the vector that has the greatest projection along the negative gradient, which amounts to solving a linear optimization problem over the domain. When \mathcal{F} is the convex

hull $conv(\mathcal{H})$ of a simpler class \mathcal{H} , this can be done using a sufficiently powerful agnostic learning primitive for \mathcal{H} . Thus we can "boost" such a primitive in a black-box manner to minimize L(f).

Let $\mathcal{H} \subset L^2(\mathbb{R}^n, D)$ be a base hypothesis class for which we have an agnostic learner with the following guarantee:

Assumption 3.1. There is an SQ learner for \mathcal{H} with the following guarantee. Let \mathcal{D} be any labeled distribution on $\mathbb{R}^n \times \mathbb{R}$ such that the marginal on \mathbb{R}^n is $D = \mathcal{N}(0, I_n)$. Given only SQ access to \mathcal{D} , the learner outputs a function $f \in \mathcal{B}(\operatorname{diam}(\mathcal{H})/2)$ such that

$$\underset{(x,y) \sim \mathcal{D}}{\mathbb{E}}[f(x)y] \ge \max_{h \in \mathcal{H}} \underset{(x,y) \sim \mathcal{D}}{\mathbb{E}}[h(x)y] - \epsilon$$

using $q(n, \epsilon, \tau)$ queries of tolerance τ .

Notice that we do not require f to lie in \mathcal{H} , i.e. the learner is allowed to be improper, but we do require it to have norm at most $\operatorname{diam}(\mathcal{H})/2$. This is to make the competitive guarantee against \mathcal{H} meaningful, since otherwise the correlation can be made to scale arbitrarily with the norm.

With such an \mathcal{H} in place, we define $\mathcal{F} = \operatorname{conv}(\mathcal{H})$. We assume that $f^* \in \mathcal{F}$. Our objective will be to agnostically learn \mathcal{F} : to solve $\min_{f \in \mathcal{F}} L(f)$ in such a way that $L(f) - L(f^*) \leq \epsilon$. To be able to use Frank-Wolfe, we require some assumptions on the loss function ℓ .

Assumption 3.2. The loss function $\ell : \mathbb{R} \times \mathbb{R} \to \mathbb{R}$ is β -smoothly convex in its first parameter.

From this assumption, orresponding properties of the loss functional L now follow. First we establish the subgradient, which will itself be an element of $L^2(\mathbb{R}^n, D)$, i.e. a function from \mathbb{R}^n to \mathbb{R} . Let $f, h : \mathbb{R}^n \to \mathbb{R}$. Observe that at for every $x \in \mathbb{R}^n$, $y \in \mathbb{R}$, the subgradient property of ℓ tells us that

$$\ell(f(x) + h(x), y) - \ell(f(x), y) \ge \partial_1 \ell(f(x), y) h(x).$$

Taking expectations over $(x, y) \sim D_{f^*}$, this yields

$$\begin{split} L(f+h) - L(f) &\geq \mathop{\mathbb{E}}_{(x,y) \sim D_{f^*}} [\partial_1 \ell(f(x),y) h(x)] \\ &= \mathop{\mathbb{E}}_{x \sim D} [\mathop{\mathbb{E}}_{y|x} [\partial_1 \ell(f(x),y)] h(x)] \\ &= \langle s,h \rangle, \end{split}$$

where

$$s: x \mapsto \mathbb{E}_{y|x}[\partial_1 \ell(f(x), y)] = \mathbb{E}_{y \sim f^*(x)}[\partial_1 \ell(f(x), y)]$$

is thus a subgradient of L at f. β -smooth convexity is also easily established. Taking expectations over $(x, y) \sim D_{f^*}$ of the inequality

$$\ell(f(x) + h(x), y) - \ell(f(x), y) - \partial_1 \ell(f(x), y) h(x) \le \frac{\beta}{2} h(x)^2,$$

we get

$$L(f+h) - L(f) - \langle s, h \rangle \le \frac{\beta}{2} ||h||^2$$

for the same subgradient s. By smooth convexity, this subgradient is unique and so we can say that the gradient of L at f is given by $\nabla L(f): x \mapsto \mathbb{E}_{y \sim f^*(x)}[\partial_1 \ell(f(x), y)].$

Example 3.3. The canonical example is the squared loss functional, with $\ell_{sq}(a,b) = (a-b)^2$, which is 2-smoothly convex. Here the gradient has a very simple form, since $\partial_1 \ell_{sq}(a,b) = 2(a-b)$, and so

$$\underset{y \sim f^*(x)}{\mathbb{E}} [\partial_1 \ell_{\text{sq}}(f(x), y)] = \underset{y \sim f^*(x)}{\mathbb{E}} [2(f(x) - y)] = 2(f(x) - f^*(x)),$$

i.e. $\nabla L_{sq}(f) = 2(f - f^*)$. In fact, it is easily calculated that

$$\begin{split} L_{\text{sq}}(f) &= \underset{(x,y) \sim D_{f^*}}{\mathbb{E}}[(f(x) - y)^2] = \underset{(x,y) \sim D_{f^*}}{\mathbb{E}}[f(x)^2] - 2 \underset{(x,y) \sim D_{f^*}}{\mathbb{E}}[f(x)y] + \underset{(x,y) \sim D_{f^*}}{\mathbb{E}}[y^2] \\ &= \underset{x \sim D}{\mathbb{E}}[f(x)^2] - 2 \underset{x \sim D}{\mathbb{E}}[f(x) \, \mathbb{E}[y|x]] + \underset{(x,y) \sim D_{f^*}}{\mathbb{E}}[y^2] \\ &= \|f\|^2 - 2\langle f, f^* \rangle + 1, \end{split}$$

It is also useful to note that

$$L_{sq}(f) - L_{sq}(f^*) = ||f - f^*||^2.$$
(2)

Frank-Wolfe using statistical queries We see that our loss functional is a β -smoothly convex functional on the space $L^2(\mathbb{R}^n, D)$. We can now use Frank-Wolfe if we can solve its main subproblem: finding an approximate solution to $\max_{h \in \mathcal{F}} \langle h, -\nabla L(f) \rangle$, where f is the current hypothesis during some iteration. Since this is a linear optimization objective and $\mathcal{F} = \text{conv}(\mathcal{H})$, this is the same as solving $\max_{h \in \mathcal{H}} \langle h, -\nabla L(f) \rangle$. This is almost the guarantee that Assumption 3.1 gives us, but some care is in order. What we have SQ access to is the labeled distribution D_{f^*} on $\mathbb{R}^n \times \{\pm 1\}$. It is not clear that we can rewrite the optimization objective in such a way that

$$\max_{h \in \mathcal{H}} \underset{x \sim D}{\mathbb{E}} [-h(x)\nabla L(f)(x)] = \max_{h \in \mathcal{H}} \underset{(x,y') \sim \mathcal{D}}{\mathbb{E}} [h(x)y']$$
 (3)

for some distribution \mathcal{D} on $\mathbb{R}^n \times \mathbb{R}$ that we can simulate SQ access to. Naively, we might try to do this by letting \mathcal{D} be the distribution of $(x, -\nabla L(f)(x))$ for $x \sim D$, so that a query $\phi: \mathbb{R} \times \mathbb{R} \to \mathbb{R}$ to \mathcal{D} can be answered with $\mathbb{E}_{(x,y')\sim\mathcal{D}}[\phi(x,y')] = \mathbb{E}_{x\sim D}[\phi(x,-\nabla L(f)(x))]$. But the issue is that in general $\nabla L(f)(x)$ will depend on $f^*(x)$, which we do not know—all we have access to is D_{f^*} .

It turns out that for the loss functions we are interested in, we can indeed find a suitable such \mathcal{D} . We turn to the details now.

4 Functional gradient descent guarantees on surrogate loss

The functional GD approach applied directly to squared loss would allow us to learn $\mathcal{F} = \text{conv}(\mathcal{H})$ using a learner for \mathcal{H} (that satisfied Assumption 3.1). But by considering a certain

surrogate loss, we can use the same learner to actually learn $\psi \circ \mathcal{F} = \{\psi \circ f \mid f \in \mathcal{F}\}$ for an outer activation function ψ . This is particularly useful as we can now capture p-concepts corresponding to functions in \mathcal{F} by using a suitable $\psi : \mathbb{R} \to [-1,1]$. For example, the common softmax activation corresponds to taking $\psi = \tanh$.

Assume that $\mathbb{E}[y|x] = \psi(f^*(x))$ for some activation $\psi : \mathbb{R} \to \mathbb{R}$ which is non-decreasing and λ -Lipschitz. Instead of the squared loss, we will consider the following surrogate loss:

$$\ell_{\mathsf{sur}}(a,b) = \int_0^a (\psi(u) - b) du.$$

It is not hard to see that $\ell_{\mathsf{sur}}(a,b)$ is convex in its first parameter due to the non-decreasing property of ψ , and that $\partial_1 \ell_{\mathsf{sur}}(a,b) = \psi(a) - b$. In fact it is λ -smoothly convex:

$$\begin{split} \ell_{\mathsf{sur}}(a+t,b) &- \ell_{\mathsf{sur}}(a,b) - \partial_1 \ell_{\mathsf{sur}}(a,b)t \\ &= \int_0^{a+t} (\psi(u) - b) du - \int_0^a (\psi(u) - b) du - (\psi(a) - b)t \\ &= \int_a^{a+t} (\psi(u) - b) du - (\psi(a) - b)t \\ &= \int_a^{a+t} (\psi(u) - \psi(a)) du \\ &\leq \int_a^{a+t} \lambda(u-a) du \\ &= \frac{\lambda t^2}{2}. \end{split}$$

The gradient of the surrogate loss functional, $L_{\mathsf{sur}}(f) = \mathbb{E}_{(x,y) \sim D_{\psi \circ f^*}}[\ell_{\mathsf{sur}}(f(x),y)]$, is given by

$$\nabla L_{\operatorname{sur}}(f): x \mapsto \underset{y \sim \psi(f^*(x))}{\mathbb{E}} [\partial_1 \ell_{\operatorname{sur}}(f(x), y)] = \psi(f(x)) - \psi(f^*(x)),$$

i.e. $\nabla L_{\mathsf{sur}}(f) = \psi \circ f - \psi \circ f^*$.

We still need to show that the Frank–Wolfe subproblem can be solved using access to just $D_{\psi \circ f^*}$. Observe that

$$\begin{split} \underset{x \sim D}{\mathbb{E}} [-h(x) \nabla L_{\text{sur}}(f)(x)] &= \underset{x \sim D}{\mathbb{E}} \left[h(x) (\psi(f^*(x)) - \psi(f(x))) \right] \\ &= \underset{x \sim D}{\mathbb{E}} \left[h(x) \left(\underset{y \sim \psi(f^*(x))}{\mathbb{E}} [y] - \psi(f(x)) \right) \right] \\ &= \underset{(x,y) \sim D_{\psi \circ f^*}}{\mathbb{E}} [h(x) (y - \psi(f(x)))] \\ &= \underset{(x,y') \sim D}{\mathbb{E}} [h(x) y'], \end{split}$$

where \mathcal{D} is the distribution of $(x, y - \psi(f(x)))$ for $(x, y) \sim D_{\psi \circ f^*}$. We can easily simulate SQ access to this using $D_{\psi \circ f^*}$: if ϕ is any query to \mathcal{D} , then

$$\mathbb{E}_{(x,y')\sim\mathcal{D}}[\phi(x,y')] = \mathbb{E}_{(x,y)\sim D_{\psi\circ f^*}}[\phi(x,y-\psi(f(x)))] = \mathbb{E}_{(x,y)\sim D_{\psi\circ f^*}}[\phi'(x,y)]$$
(4)

for the modified query $\phi'(x,y) = \phi(x,y-\psi(f(x)))$. This means we can rewrite the optimization objective to fit the form in Eq. (3). Thus for our surrogate loss, Assumption 3.1 allows us to solve the Frank-Wolfe subproblem, giving us Algorithm 2 for learning \mathcal{F} .

Algorithm 2 Frank-Wolfe for solving $\min_{f \in \mathcal{F}} L_{\mathsf{sur}}(f)$

Start with an arbitrary $f_0 \in \mathcal{B}(\operatorname{diam}(\mathcal{H})/2)$.

for $t = 0, \dots, T$ do

Let γ_t be $\frac{2}{t+2}$.

Let \mathcal{D}_t be the distribution of $(x, y - \psi(f_t(x)))$ for $(x, y) \sim D_{\psi \circ f^*}$.

Using Assumption 3.1, find $h \in \mathcal{B}(\operatorname{diam}(\mathcal{H})/2)$ such that

$$\mathbb{E}_{(x,y')\sim\mathcal{D}_t}[h(x)y'] \ge \max_{h'\in\mathcal{H}} \mathbb{E}_{(x,y')\sim\mathcal{D}_t}[h'(x)y'] - \frac{1}{2}\gamma_t\lambda\operatorname{diam}(\mathcal{H})^2$$

Let $f_{t+1} = (1 - \gamma_t) f_t + \gamma_t h$.

end for

Theorem 4.1. Let \mathcal{H} be a class for which Assumption 3.1 holds, and let $\mathcal{F} = \operatorname{conv}(\mathcal{H})$. Given SQ access to $D_{\psi \circ f^*}$ for a known non-decreasing λ -Lipschitz activation ψ and an unknown $f^* \in \mathcal{F}$, suppose we wish to learn $\psi \circ f^*$ in terms of surrogate loss, i.e. to minimize $L_{\mathsf{sur}}(f)$. Then after T iterations of Algorithm 2, we have the following guarantee:

$$L_{\mathsf{sur}}(f_T) - L_{\mathsf{sur}}(f^*) \le \frac{4\lambda \operatorname{diam}(\mathcal{H})^2}{T+2}.$$

In particular, we can achieve $L_{\text{sur}}(f_T) - L_{\text{sur}}(f^*) \leq \epsilon$ after $T = O(\frac{\lambda \operatorname{diam}(\mathcal{H})^2}{\epsilon})$ iterations. Assuming our queries are of tolerance τ , the total number of queries used is at most $Tq(n, \epsilon/4, \tau) = O(\frac{\lambda \operatorname{diam}(\mathcal{H})^2}{\epsilon}q(n, \epsilon/4, \tau))$.

Proof. By the preceding discussion, the surrogate loss functional is λ -smoothly convex, and Algorithm 2 is a valid special case of Algorithm 1, with $\mathcal{Z} = \mathcal{B}(\operatorname{diam}(\mathcal{H})/2)$ and $\mathcal{Z}' = \operatorname{conv}(\mathcal{F})$. Thus the guarantee follows directly from Theorem 2.5 (setting $\delta = 1$).

To bound the number of queries, observe that it is sufficient to run for $T = \frac{4\lambda \operatorname{diam}(\mathcal{H})^2}{\epsilon} - 2$ rounds. In the t^{th} iteration, we invoke Assumption 3.1 with

$$\epsilon' = \frac{1}{2} \gamma_t \lambda \operatorname{diam}(\mathcal{H})^2 = \frac{\lambda \operatorname{diam}(\mathcal{H})^2}{t+2} \ge \frac{\lambda \operatorname{diam}(\mathcal{H})^2}{T+2} = \frac{\epsilon}{4}.$$

Since $q(n, \epsilon', \tau) \leq q(n, \epsilon/4, \tau)$, the bound follows.

Lastly, we can show that minimizing surrogate loss also minimizes the squared loss. Observe first that $\nabla L_{\mathsf{sur}}(f^*) = 0$. Thus, applying Eq. (1) with $z = f^*$ and z' = f, we obtain

$$L_{\operatorname{sur}}(f) - L_{\operatorname{sur}}(f^{*}) \geq \frac{1}{2\lambda} \|\nabla L_{\operatorname{sur}}(f) - \nabla L_{\operatorname{sur}}(f^{*})\|^{2}$$

$$= \frac{1}{2\lambda} \|\psi \circ f - \psi \circ f^{*}\|^{2}$$

$$= \frac{1}{2\lambda} (L_{\operatorname{sq}}(\psi \circ f) - L_{\operatorname{sq}}(\psi \circ f^{*})),$$
(5)

where L_{sq} is squared loss w.r.t. $D_{\psi \circ f^*}$ and the last equality is Eq. (2). In particular, Eq. (5) implies that $\psi \circ f$ achieves the following L^2 error with respect to $\psi \circ f^*$:

$$\|\psi \circ f - \psi \circ f^*\| \le \sqrt{2\lambda \left(L_{\mathsf{sur}}(f) - L_{\mathsf{sur}}(f^*)\right)}. \tag{6}$$

5 Lower bounds on learning ReLUs, sigmoids, and halfspaces

The machinery so far has shown that if we could agnostically learn a single unit (e.g. a ReLU or a sigmoid), we could learn depth-two neural networks composed of such units. Since we have lower bounds on the latter problem, this yields the following lower bounds on the former.

Theorem 5.1. Let $\mathcal{H}_{ReLU} = \{x \mapsto \pm \text{ReLU}(w \cdot x) \mid ||w||_2 \leq 1\}$ be the class of ReLUs on \mathbb{R}^n with unit weight vectors. Suppose that Assumption 3.1 holds for \mathcal{H}_{ReLU} . Then for any ϵ , there exists $\tau = n^{-\Theta(\epsilon^{-1/12})}$ such that $q(n, \epsilon, \tau) \geq 2^{n^c} \epsilon$ for some 0 < c < 1/2.

Proof. Since all our lower bound proofs are similar, to set a template we lay out all the steps as clearly as possible.

- Consider the class \mathcal{G} from Theorem 2.3 instantiated with $\psi = \tanh$ (which is 1-Lipschitz, so $\lambda = 1$) and $\phi = \text{ReLU}$. By the conditions on the weights, we see that $\mathcal{G} \subseteq \tanh \circ \mathcal{F}_{\text{ReLU}}$, where $\mathcal{F}_{\text{ReLU}} = \text{conv}(\mathcal{H}_{\text{ReLU}})$. This construction has a free parameter k, which we will set based on ϵ .
- By our main reduction (Assumption 3.1 and Theorem 4.1), we can learn $\tanh \circ \mathcal{F}_{ReLU}$ with respect to L_{sur} up to agnostic error ϵ using $O(\frac{1}{\epsilon}q(n,\frac{\epsilon}{4},\tau))$ queries of tolerance τ . By Eq. (6), this implies learning \mathcal{G} up to L^2 error $\sqrt{2\epsilon}$.
- We know that learning \mathcal{G} should be hard. Specifically, Corollary 2.4(a) states that if $\epsilon' = \Theta(1/k^6)$ and the queries are of tolerance $\tau = n^{-\Theta(k)}$, then learning up to L^2 error ϵ' should require 2^{n^c} queries.
- The loss our reduction achieves is $\epsilon' = \sqrt{2\epsilon}$, so we require $\sqrt{2\epsilon} \leq \Theta(1/k^6)$ for the bound to hold. Accordingly, we pick $k = \Theta(\epsilon^{-1/12})$, so that $\tau = n^{-\Theta(k)} = n^{-\Theta(\epsilon^{-1/12})}$.
- Thus we must have $\frac{1}{\epsilon}q(n,\frac{\epsilon}{4},\tau)\geq 2^{n^c}$. Rearranging and rescaling ϵ gives the result.

Theorem 5.2. Let $\mathcal{H}_{\sigma} = \{x \mapsto \pm \sigma(w \cdot x) \mid ||w||_2 \leq 1\}$, where σ is the standard sigmoid, be the class of sigmoid units on \mathbb{R}^n with unit weight vectors. Suppose that Assumption 3.1 holds for \mathcal{H}_{σ} . Then for any ϵ , there exists $\tau = n^{-\Theta((\log 1/\epsilon)^2)}$ such that $q(n, \epsilon, \tau) \geq 2^{n^c} \epsilon$ for some 0 < c < 1/2.

Proof. Very similar to the above. We instantiate \mathcal{G} with $\psi = \tanh$, $\phi = \sigma$, and observe that $\mathcal{G} \subseteq \tanh \circ \operatorname{conv}(\mathcal{H}_{\sigma})$ and that $\operatorname{diam}(\mathcal{H}_{\sigma}) \leq 2$. In this case, Corollary 2.4(b) tells us that we require $\sqrt{2\epsilon} \leq e^{-\Theta(\sqrt{k})}$ for the lower bound to hold, so we pick $k = (\log 1/\epsilon)^2$. The result now follows exactly as before.

 $[\]overline{^{1}}$ We use $\pm \overline{\text{ReLU}}$ for simplicity. Any learner can handle this by doing a bit flip on its own.

We also obtain a lower bound on the class of halfspaces. The traditional way of phrasing agnostic learning for Boolean functions is in terms of the 0-1 loss, and it is not immediately obvious that the correlation loss guarantee of Assumption 3.1 is equivalent. But in Appendix E, we show that with a little care, they are indeed effectively equivalent. Note that for Boolean functions, functional GD is not essential; existing distribution-specific boosting methods [KK09, Fel10] can also give us similar results here.

Theorem 5.3. Let $\mathcal{H}_{\mathsf{hs}} = \{x \mapsto \operatorname{sign}(w \cdot x) \mid ||w||_2 \leq 1\}$ be the class of halfspaces on \mathbb{R}^n with unit weight vectors. Suppose that Assumption 3.1 holds for $\mathcal{H}_{\mathsf{hs}}$. Then for any ϵ , there exists $\tau = n^{-\Theta(1/\epsilon)}$ such that $q(n, \epsilon, \tau) \geq 2^{n^c} \epsilon^3$ for some 0 < c < 1/2.

Proof. To approximate the sign function using a Lipschitz function, we define $\operatorname{sign}(x)$ to be -1 for $x \leq -1/k$, 1 for $x \geq 1/k$, and linearly interpolate in between. This function is (k/2)-Lipschitz. We claim that \mathcal{G} instantiated with $\psi = \phi = \operatorname{sign}$ satisfies $\mathcal{G} \subseteq \operatorname{sign} \circ \operatorname{conv}(\mathcal{H}_{hs})$, with $\operatorname{diam}(\mathcal{G}) = 2$. This is because as noted in Theorem 2.3, \mathcal{G} has weights $a_i \in \{\pm 1/k\}$, so the sum of halfspaces inside ψ is always a multiple of 1/k, and sign behaves the same as sign.

Theorem 4.1 now lets us learn \mathcal{G} up to agnostic error ϵ (and hence L^2 error $\sqrt{2k\epsilon}$, by Eq. (6)) using $O(\frac{k^2}{\epsilon}q(n,\epsilon/4,\tau))$ queries of tolerance τ . By Corollary 2.4(c), we only need $\sqrt{2k\epsilon} \leq \Theta(1)$ for the lower bound to hold, so we may take $k = \Theta(1/\epsilon)$ to get a lower bound of 2^{n^c} . Thus $\frac{k^2}{\epsilon}q(n,\epsilon/4,\tau) \geq 2^{n^c}$, and rearrangement gives the result.

6 Lower bounds on learning general non-polynomial activations

Here we extend our lower bounds to general non-polynomial activations $\phi: \mathbb{R} \to \mathbb{R}$, by which we mean functions which have an infinite Hermite series $\phi = \sum_a \widehat{\phi}_a H_a$, where the H_a are the normalized probabilists' Hermite polynomials. We will again work with the class \mathcal{G} from Theorem 2.3, instantiated with this ϕ and $\psi = \tanh$. In Appendix B, we define this construction formally, letting g be the inner function and f be $\psi \circ g$.

To apply our framework, we need a norm lower bound on f. In Lemma B.1 we show that ||g|| is determined only by k, the number of hidden units (there k=2m), and the Hermite expansion of ϕ . The reason we require an infinite Hermite series for ϕ is so that this lower bound, viewed as a function of k, is nonzero for infinitely many k. This then implies that $f = \tanh \circ g$ must be nonzero for infinitely many k. Its norm can only possibly be a function of ϕ and k. In particular, we may assume that it satisfies a norm lower bound $||f|| \geq \beta(k)$, where β is a function only of k that is nonzero for infinitely many k. Here we view the dependence on ϕ as constant.

A few remarks are in order as to how such a bound $\beta(k)$ may be quantitatively established. If ϕ is either bounded or exhibits only polynomial growth, then the bound on ||g|| (Lemma B.1) gives a corresponding lower bound on ||f|| that is also purely a function of k. If ϕ is bounded, the calculation is straightforward and very similar to the $\phi = \sigma$ case (Lemma B.6). If ϕ grows only like a polynomial, then one can use a truncation argument similar to the $\phi = \text{ReLU}$ case (Lemma B.4).

By Theorem 2.2 and Corollary 2.4, our lower bound of 2^{n^c} on learning \mathcal{G} holds for $\epsilon \leq \beta(k)/3$. Since we can pick k as we like, let us say that for all sufficiently small ϵ , we

can achieve $\epsilon \leq \beta(k)/3$ by taking $k = k(\epsilon) = 3\beta^{-1}(\epsilon)$. The corresponding tolerance is then $\tau = n^{-\Theta(k(\epsilon))}$, which is still inverse superpolynomial in n.

We now get the following lower bound on learning $\mathcal{H} = \{x \mapsto \phi(w \cdot x) \mid ||w||_2 \leq 1\}$, again by the same arguments as in Section 5. We assume that $||\phi|| \leq R$ for some R, so that $\operatorname{diam}(\mathcal{H}) \leq 2R$.

Theorem 6.1. Suppose that Assumption 3.1 holds for \mathcal{H} . Then for all sufficiently small ϵ and $\tau = n^{-\Theta(k(\epsilon))}$, $q(n, \epsilon, \tau) \geq 2^{n^c} \frac{\epsilon}{R^2}$ for some 0 < c < 1/2.

Proof. We have $\mathcal{G} \subseteq \tanh \circ \operatorname{conv}(\mathcal{H})$. By functional GD wrt surrogate loss (Theorem 4.1), we see that we can learn \mathcal{G} up to L^2 error $\sqrt{2\epsilon}$ using $O(\frac{R^2}{\epsilon}q(n,\epsilon,\tau))$ queries of tolerance τ , but we must have $O(\frac{R^2}{\epsilon}q(n,\epsilon,\tau)) \leq 2^{n^c}$.

7 Lower bounds on learning monomials

In this section we show lower bounds against agnostically learning monomials with respect to the Gaussian, establishing Theorem 1.4. Let \mathcal{H}_{mon} be the class of all multilinear monomials of total degree d on \mathbb{R}^n . Clearly $|\mathcal{H}_{mon}| = \binom{n}{d} = n^{\Theta(d)}$. For any two distinct multilinear monomials f, g, clearly $\langle f, g \rangle = 0$ and moreover $\langle \tanh \circ f, \tanh \circ g \rangle = 0$ as well. Thus the class $\mathcal{G} = \tanh \circ \mathcal{H}_{mon}$ consists entirely of orthogonal functions. By [GGJ⁺20, Lemma 2.6], SDA $(\mathcal{G}, \gamma) \geq |\mathcal{G}| \gamma = n^{-\Theta(d)} \gamma$.

We still need a norm lower bound on \mathcal{G} .

Lemma 7.1. Let $x_S = \prod_{i \in S} x_i$ be an arbitrary degree-d multilinear monomial on \mathbb{R}^n , where $S \subseteq [n]$ is a subset of size d. Then $\|\tanh \circ x_S\| \ge \exp(-\Theta(d))$.

Proof. Observe first that $||x_S|| = 1$. By Paley-Zygmund, we have

$$\mathbb{P}[x_S^2 \ge \theta \, \mathbb{E}[x_S^2]] \ge (1 - \theta)^2 \frac{\mathbb{E}[x_S^2]^2}{\mathbb{E}[x_S^4]}.$$

By picking $\theta = 1/2$, say, and using the fact that by Gaussian hypercontractivity,

$$\frac{\mathbb{E}[x_S^2]^2}{\mathbb{E}[x_S^4]} = \prod_{i \in S} \frac{\mathbb{E}[x_i^2]^2}{\mathbb{E}[x_i^4]} \ge \exp(-\Theta(d)),$$

we get that $\mathbb{P}[|x_S| \ge 1/2] \ge \exp(-\Theta(d))$.

Now since tanh is monotonic and odd, we have

$$\mathbb{E}[\tanh(x_S)^2] \ge \tanh(1/2)^2 \,\mathbb{P}[|x_S| \ge 1/2] \ge \exp(-\Theta(d)).$$

By Theorem 2.2 with $\beta = \exp(-\Theta(d))$, we get that for any $\epsilon \leq \exp(-\Theta(d))$ and using queries of tolerance $\tau \leq \epsilon^2$, learning \mathcal{G} up to L^2 error ϵ takes at least $\mathrm{SDA}(\mathcal{G}, \tau^2) \geq n^{\Theta(d)} \tau^2$ queries.

Now we can use the same arguments as in Section 5 to prove the following.

Theorem 7.2. Suppose that Assumption 3.1 holds for \mathcal{H}_{mon} . Then for any $\epsilon \leq \exp(-\Theta(d))$ and $\tau \leq \epsilon^2$, $q(n, \epsilon, \tau) \geq n^{\Theta(d)} \tau^{5/2}$.

Proof. Observe that $\mathcal{G} \subseteq \tanh \circ \operatorname{conv}(\mathcal{H}_{\mathsf{mon}})$, and $\operatorname{diam}(\mathcal{H}_{\mathsf{mon}}) \leq 2$. Using the surrogate loss with $\psi = \tanh$, Assumption 3.1 and Theorem 4.1 tell us that we can learn $\tanh \circ \operatorname{conv}(\mathcal{H}_{\mathsf{mon}})$ up to L^2 error $\sqrt{2\epsilon}$ (again by Eq. (2)) in $O(\frac{1}{\epsilon}q(n,\epsilon,\tau))$ queries of tolerance τ . By our lower bound for \mathcal{G} , we must have $\frac{1}{\epsilon}q(n,\epsilon,\tau) \geq n^{\Theta(d)}\tau^2$, or $q(n,\epsilon,\tau) \geq n^{\Theta(d)}\tau^{5/2}$ (since $\epsilon \geq \sqrt{\tau}$). \square

8 Upper bounds on learning ReLUs and sigmoids

We use a variant of the classic low-degree algorithm ([LMN93]; see also [KKMS08]) to provide simple upper bounds for agnostically learning ReLUs and sigmoids. With respect to $D = \mathcal{N}(0, I_n)$, the δ -approximate degree of a function $f : \mathbb{R}^n \to \mathbb{R}$ is the smallest d such that there exists a degree-d polynomial p satisfying $||f - p|| \le \delta$. We show that for any class of δ -approximate degree d, picking $\delta = O(\epsilon)$ and simply estimating the Hermite coefficients of $x \mapsto \mathbb{E}[y|x]$ up to degree d yields an agnostic learner up to error ϵ , one that satisfies Assumption 3.1. We assume bounded labels, say $y \in [-C, C]$ for some constant C.

Let \mathcal{D} be a distribution on $\mathbb{R}^n \times \mathbb{R}$ such that the marginal on \mathbb{R}^n is $\mathcal{N}(0, I_n)$. Let $f_{\mathsf{cmf}}(x) = \mathbb{E}[y|x]$ denote the conditional mean function of \mathcal{D} , and note that $||f_{\mathsf{cmf}}|| \leq C$. Observe that for any f, the correlation $\mathbb{E}_{(x,y)\sim\mathcal{D}}[f(x)y]$ equals $\langle f, f_{\mathsf{cmf}} \rangle$. Let \mathcal{H} be a hypothesis class with δ -approximate degree d (δ to be determined), and let $R = \dim(\mathcal{H})/2$. Let $h_{\mathsf{opt}} \in \mathcal{H}$ achieve $\max_{h \in \mathcal{H}} \langle h, f_{\mathsf{cmf}} \rangle$.

Our algorithm will be based on approximating the low-degree Hermite coefficients of f_{cmf} , which is equivalent to performing polynomial L^2 regression. It is well-known that in this context, where d is the δ -approximate degree, polynomial L^1 regression up to degree d gives a squared loss guarantee of δ [KKMS08]. But we will not be able to use this result directly since what we seek is a correlation guarantee. Instead, our approach will involve a sequence of inequalities relating the correlation achieved by f_{cmf} , h_{opt} , and their degree-d approximations. A slight subtlety to keep in mind is that correlation can always be increased by scaling the function. This means that wherever scaling is possible, we have to take some care to rescale functions to have the maximum allowed norm, R.

Let $h_{\mathsf{opt}}^{\leq d}$ and $f_{\mathsf{cmf}}^{\leq d}$ be the Hermite components of degree at most d of h_{opt} and f_{cmf} respectively. Let $\tilde{f}_{\mathsf{cmf}}^{\leq d} = \frac{R}{\|f_{\mathsf{cmf}}^{\leq d}\|} f_{\mathsf{cmf}}^{\leq d}$. Among polynomials of degree d in $\mathcal{B}(R)$, it is easy to see that $\tilde{f}_{\mathsf{cmf}}^{\leq d}$ maximizes $\langle f, f_{\mathsf{cmf}} \rangle$, so that

$$\langle \tilde{f}_{\mathsf{cmf}}^{\leq d}, f_{\mathsf{cmf}} \rangle \geq \langle h_{\mathsf{opt}}^{\leq d}, f_{\mathsf{cmf}} \rangle.$$

Our agnostic learner will look to approximate $\tilde{f}_{\mathsf{cmf}}^{\leq d}$ by outputting p defined as follows. Suppose $f_{\mathsf{cmf}} = \sum_{I \in \mathbb{N}^n} \alpha_I H_I$, where H_I is the multivariate Hermite polynomial of index I. For each I of total degree at most d, which we denote as $|I| \leq d$, let β_I be our estimate of $\alpha_I = \langle f_{\mathsf{cmf}}, H_I \rangle$ to within tolerance τ (to be determined). This can be done using $n^{O(d)}$

queries of tolerance τ . Let $\tilde{f} = \sum_{|I| \leq d} \beta_I H_I$, and finally let $p = \frac{R}{\|\tilde{f}\|} \tilde{f}$. We have

$$\begin{split} \|\tilde{f}_{\mathsf{cmf}}^{\leq d} - p\|^2 &= R^2 \left\| \frac{f_{\mathsf{cmf}}^{\leq d}}{\|f_{\mathsf{cmf}}^{\leq d}\|} - \frac{\tilde{f}}{\|\tilde{f}\|} \right\|^2 \\ &= R^2 \left\| \frac{f_{\mathsf{cmf}}^{\leq d} - \tilde{f}}{\|f_{\mathsf{cmf}}^{\leq d}\|} + \tilde{f} \left(\frac{1}{\|f_{\mathsf{cmf}}^{\leq d}\|} - \frac{1}{\|\tilde{f}\|} \right) \right\|^2 \\ &\leq 2R^2 \left(\frac{\|f_{\mathsf{cmf}}^{\leq d} - \tilde{f}\|^2}{\|f_{\mathsf{cmf}}^{\leq d}\|^2} + \|\tilde{f}\|^2 \left(\frac{1}{\|f_{\mathsf{cmf}}^{\leq d}\|} - \frac{1}{\|\tilde{f}\|} \right)^2 \right) \\ &= 2R^2 \left(\frac{\|f_{\mathsf{cmf}}^{\leq d} - \tilde{f}\|^2}{\|f_{\mathsf{cmf}}^{\leq d}\|^2} + \left(\frac{\|f_{\mathsf{cmf}}^{\leq d}\| - \|\tilde{f}\|}{\|f_{\mathsf{cmf}}^{\leq d}\|} \right)^2 \right) \\ &\leq 4R^2 \frac{\|f_{\mathsf{cmf}}^{\leq d} - \tilde{f}\|^2}{\|f_{\mathsf{cmf}}^{\leq d}\|^2} \\ &\leq \frac{4R^2 n^d \tau^2}{\|f_{\mathsf{cmf}}^{\leq d}\|^2}, \end{split} \tag{triangle ineq.}$$

since $\|\tilde{f} - f_{\mathsf{cmf}}^{\leq d}\| \leq n^{d/2} \tau$.

We claim that we can assume WLOG that $\|\tilde{f}_{\mathsf{cmf}}^{\leq d}\| \geq \epsilon/(2R)$. Indeed, we know $\max_{h \in \mathcal{H}} \langle h, f_{\mathsf{cmf}} \rangle = \langle h_{\mathsf{opt}}, f_{\mathsf{cmf}} \rangle$ and also $\|h_{\mathsf{opt}} - h_{\mathsf{opt}}^{\leq d}\| \leq \delta$. This implies that

$$\|R\|\tilde{f}_{\mathsf{cmf}}^{\leq d}\| = \langle \tilde{f}_{\mathsf{cmf}}^{\leq d}, f_{\mathsf{cmf}} \rangle \geq \langle h_{\mathsf{opt}}^{\leq d}, f_{\mathsf{cmf}} \rangle \geq \langle h_{\mathsf{opt}}, f_{\mathsf{cmf}} \rangle - C\delta_{\mathsf{opt}}$$

where the last inequality is Cauchy–Schwarz. If $\langle h_{\sf opt}, f_{\sf cmf} \rangle \leq \epsilon$ then 0 is a valid agnostic learner. Therefore, we can assume that $\langle h_{\sf opt}, f_{\sf cmf} \rangle \geq \epsilon$. Choosing $\delta = \frac{\epsilon}{2C}$, this means $\|\tilde{f}_{\sf cmf}^{\leq d}\| \geq \epsilon/(2R)$.

By Eq. (7), we then have

$$\|\tilde{f}_{\mathsf{cmf}}^{\leq d} - p\| \leq \frac{4Rn^{d/2}\tau}{\epsilon}.\tag{8}$$

Now observe that

$$\begin{split} \langle p, f_{\mathsf{cmf}} \rangle &= \langle \tilde{f}_{\mathsf{cmf}}^{\leq d}, f_{\mathsf{cmf}} \rangle + \langle p - \tilde{f}_{\mathsf{cmf}}^{\leq d}, f_{\mathsf{cmf}} \rangle \\ &\geq \langle h_{\mathsf{opt}}^{\leq d}, f_{\mathsf{cmf}} \rangle - \frac{4RCn^{d/2}\tau}{\epsilon} \qquad \qquad \text{(Eq. (8) and Cauchy–Schwarz)} \\ &= \langle h_{\mathsf{opt}}, f_{\mathsf{cmf}} \rangle + \langle h_{\mathsf{opt}}^{\leq d} - h_{\mathsf{opt}}, f_{\mathsf{cmf}} \rangle - \frac{4RCn^{d/2}\tau}{\epsilon} \\ &\geq \langle h_{\mathsf{opt}}, f_{\mathsf{cmf}} \rangle - \frac{\epsilon}{2} - \frac{4RCn^{d/2}\tau}{\epsilon}. \qquad \text{(Cauchy–Schwarz, and using } \delta C = \epsilon/2) \end{split}$$

Setting $\tau = \frac{\epsilon^2}{8RCn^{d/2}}$ gives us the desired result, namely that $\langle p, f_{\sf cmf} \rangle \geq \langle h_{\sf opt}, f_{\sf cmf} \rangle - \epsilon$. Thus we have the following theorem.

Theorem 8.1. The class \mathcal{H}_{ReLU} can be agnostically learned up to correlation ϵ (in the sense of Assumption 3.1) using $n^{O(\epsilon^{-4/3})}$ queries of tolerance $n^{-\Theta(\epsilon^{-4/3})}\epsilon$. Similarly, \mathcal{H}_{σ} can be

learned using $n^{\tilde{O}(\log^2 1/\epsilon)}$ queries of tolerance $n^{-\tilde{\Theta}(\log^2 1/\epsilon)}\epsilon^2$.

Proof. Approximating the Hermite coefficients of degree at most d takes $n^{O(d)}$ queries of tolerance $n^{-\Theta(d)}\epsilon$. As we show in Appendix C, the δ -approximate degree of unit-weight ReLUs is $O((1/\delta)^{4/3})$ and for unit-weight sigmoids it is $\tilde{O}(\log^2 1/\delta)$. The guarantees follow by the argument in the preceding discussion.

We note that our lower bounds for ReLUs and sigmoids were for queries of tolerance $n^{-\Theta(\epsilon^{-1/12})}$ and $n^{-\Theta(\log^2 1/\epsilon)}$ respectively, which nearly matches these upper bounds.

Acknowledgements

We thank the anonymous NeurIPS 2020 reviewers for their feedback.

References

- [ADHV19] Alexandr Andoni, Rishabh Dudeja, Daniel Hsu, and Kiran Vodrahalli. Attribute-efficient learning of monomials over highly-correlated variables. In Algorithmic Learning Theory, pages 127–161, 2019. 1
- [APVZ14] Alexandr Andoni, Rina Panigrahy, Gregory Valiant, and Li Zhang. Learning sparse polynomial functions. In *Proceedings of the twenty-fifth annual ACM-SIAM symposium on Discrete algorithms*, pages 500–510. SIAM, 2014. 1
- [Bac17] Francis Bach. Breaking the curse of dimensionality with convex neural networks. The Journal of Machine Learning Research, 18(1):629–681, 2017. 1
- [BHKL15] Alina Beygelzimer, Elad Hazan, Satyen Kale, and Haipeng Luo. Online gradient boosting. In *Advances in neural information processing systems*, pages 2458–2466, 2015. 1
- [Boy84] John P Boyd. Asymptotic coefficients of hermite function series. *Journal of Computational Physics*, 54(3):382–410, 1984. C
- [DFGK17] Carlos M Da Fonseca, M Lawrence Glasser, and Victor Kowalenko. Basic trigonometric power sums with applications. *The Ramanujan Journal*, 42(2):401–428, 2017. B
- [DGK+20] Ilias Diakonikolas, Surbhi Goel, Sushrut Karmalkar, Adam Klivans, and Mahdi Soltanolkotabi. Approximation Schemes for ReLU Regression. In Conference on Learning Theory, 2020. To appear. 1
- [DKKZ20] Ilias Diakonikolas, Daniel Kane, Vasilis Kontonis, and Nikos Zarifis. Algorithms and SQ Lower Bounds for PAC Learning One-Hidden-Layer ReLU Networks. In *Conference on Learning Theory*, 2020. To appear. 1, 1, 1, 2, 2.3, B
- [DKN10] Ilias Diakonikolas, Daniel M Kane, and Jelani Nelson. Bounded independence fools degree-2 threshold functions. In 2010 IEEE 51st Annual Symposium on Foundations of Computer Science, pages 11–20. IEEE, 2010. 1, 1

- [DKS17] Ilias Diakonikolas, Daniel M Kane, and Alistair Stewart. Statistical query lower bounds for robust estimation of high-dimensional gaussians and gaussian mixtures. In 2017 IEEE 58th Annual Symposium on Foundations of Computer Science (FOCS), pages 73–84. IEEE, 2017. 1
- [DKZ20] Ilias Diakonikolas, Daniel Kane, and Nikos Zarifis. Near-Optimal SQ Lower Bounds for Agnostically Learning Halfspaces and ReLUs under Gaussian Marginals. In *Advances in Neural Information Processing Systems*, 2020. 1
- [FCG20] Spencer Frei, Yuan Cao, and Quanquan Gu. Agnostic learning of a single neuron with gradient descent. arXiv preprint arXiv:2005.14426, 2020. 1, 1
- [Fel10] Vitaly Feldman. Distribution-specific agnostic boosting. In Andrew Chi-Chih Yao, editor, Innovations in Computer Science ICS 2010, Tsinghua University, Beijing, China, January 5-7, 2010. Proceedings, pages 241–250. Tsinghua University Press, 2010. 1, 5
- [Fel12] Vitaly Feldman. A complete characterization of statistical query learning with applications to evolvability. *Journal of Computer and System Sciences*, 78(5):1444–1459, 2012. 1
- [Fri01] Jerome H Friedman. Greedy function approximation: a gradient boosting machine. *Annals of statistics*, pages 1189–1232, 2001. 1
- [GGJ⁺20] Surbhi Goel, Aravind Gollakota, Zhihan Jin, Sushrut Karmalkar, and Adam Klivans. Superpolynomial Lower Bounds for Learning One-Layer Neural Networks using Gradient Descent. In *International Conference on Machine Learning*, 2020. To appear. 1, 1, 2, 2.2, 7, A.1, A.2, 1, 2, B, B, B.5
- [GKK19] Surbhi Goel, Sushrut Karmalkar, and Adam Klivans. Time/Accuracy Tradeoffs for Learning a ReLU with respect to Gaussian Marginals. In *Advances in Neural Information Processing Systems*, pages 8582–8591, 2019. 1, 1, 1
- [GKKT17] Surbhi Goel, Varun Kanade, Adam Klivans, and Justin Thaler. Reliably learning the relu in polynomial time. In *Conference on Learning Theory*, pages 1004–1042, 2017. 1
- [Haz16] Elad Hazan. Introduction to online convex optimization. Foundations and Trends in Optimization, 2(3-4):157–325, 2016. 1
- [Hil40] Einar Hille. Contributions to the theory of hermitian series ii. the representation problem. Transactions of the American Mathematical Society, 47(1):80–94, 1940. C
- [Jac97] Jeffrey C. Jackson. An efficient membership-query algorithm for learning dnf with respect to the uniform distribution. J. Comput. Syst. Sci, 55(3):414–440, 1997. 1
- [Jag13] Martin Jaggi. Revisiting Frank-Wolfe: Projection-free sparse convex optimization. In Proceedings of the 30th International Conference on Machine Learning, pages 427–435, 2013. 2, D

- [KK09] Varun Kanade and Adam Kalai. Potential-based agnostic boosting. In Advances in neural information processing systems, pages 880–888, 2009. 1, 5
- [KK14] Adam Klivans and Pravesh Kothari. Embedding hard learning problems into gaussian space. In Approximation, Randomization, and Combinatorial Optimization. Algorithms and Techniques (APPROX/RANDOM), 2014. 1, 1, 1
- [KKMS08] Adam Tauman Kalai, Adam R Klivans, Yishay Mansour, and Rocco A Servedio. Agnostically learning halfspaces. SIAM Journal on Computing, 37(6):1777–1805, 2008. 1, 1, 1, 8
- [KS94] Michael J Kearns and Robert E Schapire. Efficient distribution-free learning of probabilistic concepts. Journal of Computer and System Sciences, 48(3):464– 497, 1994. 1
- [LMN93] Nathan Linial, Yishay Mansour, and Noam Nisan. Constant depth circuits, fourier transform, and learnability. *Journal of the ACM (JACM)*, 40(3):607–620, 1993.
- [LSSS14] Roi Livni, Shai Shalev-Shwartz, and Ohad Shamir. On the computational efficiency of training neural networks. In *Advances in neural information processing systems*, pages 855–863, 2014. 1
- [MBBF00] Llew Mason, Jonathan Baxter, Peter L Bartlett, and Marcus R Frean. Boosting algorithms as gradient descent. In *Advances in neural information processing systems*, pages 512–518, 2000. 1
- [PSG19] Abhishek Panigrahi, Abhishek Shetty, and Navin Goyal. Effect of activation functions on the training of overparametrized neural nets. $arXiv\ preprint\ arXiv:1908.05660,\ 2019.$
- [SF12] Robert E. Schapire and Yoav Freund. *Boosting: Foundations and Algorithms*. The MIT Press, 2012. 1
- [SVWX17] Le Song, Santosh Vempala, John Wilmes, and Bo Xie. On the complexity of learning neural networks. In *Advances in neural information processing systems*, pages 5514–5522, 2017. 1
- [VW19] Santosh Vempala and John Wilmes. Gradient descent for one-hidden-layer neural networks: Polynomial convergence and sq lower bounds. In *COLT*, volume 99, 2019. 1, 1
- [YS19] Gilad Yehudai and Ohad Shamir. On the power and limitations of random features for understanding neural networks. In *Advances in Neural Information Processing Systems*, pages 6594–6604, 2019. 1
- [YS20] Gilad Yehudai and Ohad Shamir. Learning a single neuron with gradient methods. arXiv preprint arXiv:2001.05205, 2020. 1, 1

A SQ lower bound subtleties

A.1 Relationships between parameters

When formally stating SQ lower bounds on learning p-concepts in terms of the statistical dimension, there are some subtleties to keep in mind. These have to do with the relationships between the query tolerance, the desired final error, and the norms of the functions in the class. Let us say our queries are of tolerance τ , the final desired L^2 error $||f - f^*||$ is ϵ (which corresponds to $L(f) - L(f^*) \leq \epsilon^2$; see Eq. (2)), and that the functions in \mathcal{C} satisfy $||f^*|| \geq \beta$ for all $f^* \in \mathcal{C}$. Then

- 1. We must have $\tau < \epsilon$. To see why, first note that for any query ϕ and two functions $f,g \in \mathcal{C}$, a calculation shows that $|\mathbb{E}_{D_f}[\phi] \mathbb{E}_{D_g}[\phi]| = |\langle f g, \tilde{\phi} \rangle| \leq ||f g||$, where $\tilde{\phi}(x) = (\phi(x,1) \phi(x,-1))/2$. Thus if one has a function f such that $\epsilon < ||f f^*|| < \tau$, then no query of tolerance τ can tell them apart, but f is not ϵ -close to the target f^* .
- 2. If $\epsilon \geq \beta$, a lower bound might not be possible. This is because the 0 function trivially achieves L^2 error $||0 f^*|| = ||f^*||$. Imposing $\epsilon < \beta$ is sufficient to rule this out.
- 3. We cannot arbitrarily rescale the p-concepts to increase β since the functions must remain Boolean p-concepts. Rescaling would also increase the description length of the functions.

The lower bound in Theorem 2.2 (from [GGJ⁺20]) is proved by reducing a distinguishing problem to a learning problem. For technical reasons, we end up requiring $\tau \leq \epsilon^2$, $\epsilon \leq \beta/3$ for this reduction to go through. The points above show that these requirements are essentially necessary.

A.2 The dependence of the query lower bound on the error ϵ and the tolerance τ

The relationship between our query lower bounds, the desired error ϵ , and the tolerance τ may seem a little unusual at first sight, especially the fact that the lower bounds seem to grow weaker as ϵ grows smaller. We make some clarifying remarks here.

Fundamentally, all SQ lower bounds are bounds on how many queries it takes to distinguish certain distributions from others. When discussing a concept class \mathcal{C} , the distributions in question are the labeled distributions corresponding to concepts in the class. Learning \mathcal{C} is hard exactly insofar as it allows us to distinguish different labeled distributions arising from \mathcal{C} . Many works in the SQ literature have this structure, but we will refer to [GGJ⁺20] for formal statements.

Formally, the distinguishing problem we consider ([GGJ⁺20, Definition 4.2]) is that of distinguishing the labeled distribution D_c arising from an unknown $c \in \mathcal{C}$ from the reference distribution $D_0 = D \times \text{Unif}\{\pm 1\}$, using queries of tolerance at least τ .

There are two crucial points to keep in mind here:

1. The distinguishing problem is a fundamentally information theoretic problem, and its difficulty scales only with τ . In particular, using queries of tolerance τ , we need at least SDA(\mathcal{C}, τ^2) queries. This bound increases with τ ; in fact it often scales as $|\mathcal{C}|\tau^2$ (see ([GGJ⁺20, Theorem 4.5 and Lemma 2.6]).

2. The problem of learning up \mathcal{C} to error ϵ is hard exactly insofar as it allows us to solve the distinguishing problem (see [GGJ⁺20, Lemma 4.4]).

An important consequence is that for fixed τ , the query lower bound does not technically grow as a function of the error ϵ : it applies uniformly for all ϵ small enough that it allows the learner to solve the distinguishing problem. In other words, there is a certain "threshold" ϵ_0 such that for all $\epsilon \leq \epsilon_0$, the same query lower bound holds. As noted in point (3) of the previous subsection, this threshold can be taken to be $\beta/3$, where β is such that $||c|| \geq \beta$ for all $c \in \mathcal{C}$.

But at the same time, as noted in point (1) in the previous subsection, it is necessary that $\tau < \epsilon$ (and for the reduction it suffices to have $\tau \le \epsilon^2$). If $\tau \ge \epsilon$, learning up to error ϵ is simply impossible.

With all this in mind, we can now answer the question of why our lower bounds seem to grow weaker as ϵ grows smaller: it is essentially because τ grows smaller as well, so that we get a series of incomparable (though still exponential) bounds due to the tradeoffs between query complexity, τ , and ϵ .

B Bounding the function norms of the [DKKZ20] construction

We shall consider the following slight rescaling of the functions of [DKKZ20]. For activation functions $\psi, \phi : \mathbb{R} \to \mathbb{R}$, we have $g, f : \mathbb{R}^2 \to \mathbb{R}$ defined as follows.

$$g(x) = \frac{1}{2m} \sum_{i=1}^{2m} (-1)^i \phi\left(x_1 \cos\frac{i\pi}{m} + x_2 \sin\frac{i\pi}{m}\right) = \frac{1}{2m} \sum_{i=1}^{2m} (-1)^i \phi\left(x \cdot w_i\right)$$
$$f(x) = \psi(g(x)),$$

where $w_i = (\cos \frac{i\pi}{m}, \sin \frac{i\pi}{m})$. The number of hidden units is k = 2m. We will assume that m is even.

The hard functions from $\mathbb{R}^n \to \mathbb{R}$ are then given by $f_A(x) = f(Ax)$ for certain matrices $A \in \mathbb{R}^{2 \times d}$ with $AA^T = I_2$. For $x \sim \mathcal{N}(0, I_d)$, Ax has the distribution $\mathcal{N}(0, I_2)$. So for the purposes of the norm calculation, and hence throughout this section, we will work directly with $\mathcal{N}(0, I_2)$. We will start by considering the norm of g. This can then be used to control the norm of f via arguments similar to those in $[GGJ^+20]$.

Lemma B.1. Let $g: \mathbb{R}^2 \to \mathbb{R}$ be as defined above, and assume m is even. Assume the standard Hermite expansion of ϕ is given by $\phi = \sum_a \widehat{\phi}_a H_a$, where the H_a are the normalized probabilists' Hermite polynomials. Under $\mathcal{N}(0, I_2)$,

$$||g||^2 = \Omega \left(\sum_{\substack{a \gg m \ a \ even}} \frac{\widehat{\phi}_a^2}{\sqrt{a}} \right).$$

(For practical purposes, the asymptotic behavior of this expression is captured faithfully when we begin indexing from say a=100m.)

Proof. We have

$$||g||^{2} = \mathbb{E}[g(x)^{2}] = \frac{1}{4m^{2}} \sum_{i,j=1}^{2m} (-1)^{i} (-1)^{j} \mathbb{E}[\phi(x \cdot w_{i})\phi(x \cdot w_{j})]$$

$$= \frac{1}{4m^{2}} \sum_{i,j=1}^{2m} (-1)^{i} (-1)^{j} \mathbb{E}\left[\left(\sum_{a} \widehat{\phi}_{a} H_{a}(x \cdot w_{i})\right) \left(\sum_{b} \phi_{b} H_{b}(x \cdot w_{j})\right)\right]$$

$$= \frac{1}{4m^{2}} \sum_{i,j=1}^{2m} (-1)^{i} (-1)^{j} \left(\sum_{a} \widehat{\phi}_{a}^{2} \mathbb{E}[H_{a}(x \cdot w_{i}) H_{a}(x \cdot w_{j})]\right).$$

Now because w_i, w_j are both unit vectors with $w_i \cdot w_j = \cos \frac{(i-j)\pi}{m}$, we have that $x \cdot w_i$ and $x \cdot w_j$ are both $\mathcal{N}(0,1)$ with covariance $\cos \frac{(i-j)\pi}{m}$. Thus

$$||g||^2 = \frac{1}{4m^2} \sum_{i,j=1}^{2m} (-1)^{i+j} \left(\sum_a \widehat{\phi}_a^2 \cos^a \frac{(i-j)\pi}{m} \right)$$
$$= \frac{1}{4m^2} \sum_{i,j=1}^{2m} (-1)^{i-j} \left(\sum_a \widehat{\phi}_a^2 \cos^a \frac{(i-j)\pi}{m} \right),$$

since $(-1)^{i+j} = (-1)^{i-j}$. Now, as we range over $i, j \in [2m]$, we see that i-j=0 occurs 2m times, i-j=1 occurs 2m-1 times, and more generally i-j=t occurs 2m-|t| times. Since a term with i-j=t is exactly the same as one with i-j=-t (by the evenness of cos), we can say that for $t \neq 0$, |i-j|=t occurs 2(2m-t) times. Thus the expression above can be written as

$$||g||^{2} = \frac{1}{4m^{2}} \left(2m \left(\sum_{a} \widehat{\phi}_{a}^{2} \cos^{a} 0 \right) + \sum_{t=1}^{2m-1} 2(2m-t)(-1)^{t} \left(\sum_{a} \widehat{\phi}_{a}^{2} \cos^{a} \frac{t\pi}{m} \right) \right)$$

$$= \frac{1}{4m^{2}} \sum_{a} \widehat{\phi}_{a}^{2} \left(2m + \sum_{t=1}^{2m-1} 2(2m-t)(-1)^{t} \cos^{a} \frac{t\pi}{m} \right)$$

$$= \frac{1}{4m^{2}} \sum_{a} \widehat{\phi}_{a}^{2} S(a, m), \tag{9}$$

where

$$S(a,m) = 2m + \sum_{t=1}^{2m-1} 2(2m-t)(-1)^t \cos^a \frac{t\pi}{m}.$$

Now some algebraic manipulations are in order. By rewriting the index t as 2m-t, we get that

$$S(a,m) = 2m + \sum_{t=1}^{2m-1} 2t(-1)^{2m-t} \cos^a \frac{(2m-t)\pi}{m}$$
$$= 2m + \sum_{t=1}^{2m-1} 2t(-1)^t \cos^a \frac{t\pi}{m}.$$

Adding the two expressions for S(a, m) and dividing by 2, we get

$$S(a,m) = 2m + \sum_{t=1}^{2m-1} 2m(-1)^t \cos^a \frac{t\pi}{m}$$
$$= 2m \sum_{t=0}^{2m-1} (-1)^t \cos^a \frac{t\pi}{m}$$

This sum vanishes when a and m have different parities, i.e. if a is odd (recall that we assume m is even). For even a, we have

$$S(a,m) = 4m \sum_{t=0}^{m-1} (-1)^t \cos^a \frac{t\pi}{m}.$$

This is a trigonometric power sum with known closed form expressions. In particular, Equation 3.4 from [DFGK17, §3] (after correcting a typo) tells us that

$$T(a,m) = \sum_{t=0}^{m-1} (-1)^t \cos^a \frac{t\pi}{m} = \begin{cases} 2^{1-a} m \left(\sum_{p=1}^{\lfloor a/m \rfloor} \binom{a}{a/2 - pm/2} \right) - \sum_{p=1}^{\lfloor a/2m \rfloor} \binom{a}{a/2 - pm} \right) & a \ge 2m \\ 2^{1-a} m \sum_{p=1}^{\lfloor a/m \rfloor} \binom{a}{a/2 - pm/2} & m \le a < 2m \end{cases}$$

$$= \begin{cases} 2^{1-a} m \left(\sum_{p=1}^{\lfloor a/m \rfloor} \binom{a}{a/2 - pm/2} \right) & a \ge m \\ 0 & a < m \end{cases}$$

To get a sense for the asymptotics as $a \to \infty$, we consider $a \gg m$ (say $a \ge 100m$). In this regime the sum of binomial coefficients in the sum above is seen to be $\Omega(2^a/\sqrt{a})$ (the p=1 term alone contributes roughly $\binom{a}{a/2}$), and we get that $T(a,m) = \Omega(m/\sqrt{a})$.

This means S(a,m)=0 for odd a and $S(a,m)=4mT(a,m)=\Omega(m^2/\sqrt{a})$ for large, even a. Substituting this back into Eq. (9), we get that

$$||g||^2 = \Omega \left(\sum_{\substack{a \gg m \ a \text{ even}}} \frac{\widehat{\phi}_a^2}{\sqrt{a}} \right).$$

We can now consider the special cases of $\phi=\text{ReLU}$ and $\phi=\sigma$ (the standard sigmoid) that are of interest.

Corollary B.2. Consider g instantiated with $\phi = \text{ReLU}$. Then $||g|| = \Omega(1/m)$.

Proof. The Hermite coefficients of ReLU satisfy $\widehat{\phi}_a = \Theta(a^{-5/4})$ (Lemma C.1). Thus by

Lemma B.1,

$$||g||^2 = \Omega(\sum_{\substack{a \ge 100m \ a \text{ even}}} a^{-3}) = \Omega(1/m^2).$$

Corollary B.3. Consider g instantiated with $\phi = \sigma$, the standard sigmoid. Then $||g|| = e^{-O(\sqrt{m})}$.

Proof. The Hermite coefficients of σ asymptotically satisfy $\widehat{\phi}_a \simeq e^{-C\sqrt{a}}$ [GGJ⁺20, §A.2] for some C. Thus by Lemma B.1,

$$||g||^2 = \Omega\left(\sum_{\substack{a \ge 100m \\ a \text{ even}}} \frac{e^{-\sqrt{a}}}{\sqrt{a}}\right).$$

The result then follows by the following standard integral approximation:

$$\sum_{t=N}^{\infty} \frac{e^{-\sqrt{t}}}{\sqrt{t}} \approx \int_{N}^{\infty} \frac{e^{-\sqrt{t}}}{\sqrt{t}} \ dt = 2e^{-\sqrt{N}}.$$

We can now translate these into norm lower bounds on $f = \psi \circ g$. For us it suffices to consider $\psi = \tanh : \mathbb{R} \to [-1,1]$, which is essentially the sigmoid centered at 0. The centering at 0 and the output range being [-1,1] is what is important to us, because we use f to capture the conditional mean function of a p-concept.

Lemma B.4. Consider f instantiated with $\psi = \tanh$ and $\phi = \text{ReLU}$. Then $||f|| = \Omega(1/m^6)$.

Proof. Ideally we would like to use the norm bound on g to obtain an anti-concentration inequality of the form $\mathbb{P}[|g(x)| > t]$, and then translate that into a norm lower bound for f, but this is not immediate because g is unbounded. So we introduce the function g^T , which is the same as g except with the truncated ReLU, ReLU $^T(x) = \min(T, \text{ReLU}(x))$ (T to be determined), in place of all standard ReLUs. Clearly $|g^T(x)| \leq T$ for all x. It is also easy to see by a union bound that

$$\mathbb{P}[g(x) \neq g^T(x)] \leq 2m \underset{t \sim \mathcal{N}(0,1)}{\mathbb{P}}[\mathrm{ReLU}(t) \neq \mathrm{ReLU}^T(t)] \leq 2me^{-T^2/2},$$

since each w_i is a unit vector.

Let $\text{ReLU}_w(x)$ be shorthand for $\text{ReLU}(x \cdot w)$, and similarly ReLU_w^T . Observe first that

$$||g - g^T|| = \frac{1}{2m} \left| \sum_{i=1}^{2m} (-1)^i (\text{ReLU}_{w_i} - \text{ReLU}_{w_i}^T) \right|$$

$$\leq \frac{1}{2m} \sum_{i=1}^{2m} ||\text{ReLU}_{w_i} - \text{ReLU}_{w_i}^T)||$$

$$= ||\text{ReLU} - \text{ReLU}^T||_{\mathcal{N}(0,1)}$$

$$\leq \sqrt{e^{-\frac{T^2}{2}} \left(T^2 + 1 - \frac{T}{\sqrt{2\pi}} \right)}$$

where the third equality again uses the fact the w_i are unit vectors, and the last inequality is Lemma B.5. By picking $T = \Theta(m)$, this coupled with the fact that $||g|| = \Omega(1/m)$ (Corollary B.2) tells us that $||g^T|| = \Omega(1/m)$ as well.

This bound on $||g^T||$ yields an anti-concentration inequality for g^T as follows:

$$||g^T||^2 = \mathbb{E}[g^T(x)^2] \le t^2 \mathbb{P}[|g^T(x)| \le t] + T^2 \mathbb{P}[|g^T(x)| > t] = t^2 + (T^2 - t^2) \mathbb{P}[|g^T(x)| > t],$$

so that

$$\mathbb{P}[|g^T(x)| > t] \ge \frac{||g^T||^2 - t^2}{T^2 - t^2}.$$

Recall that $\mathbb{P}[g(x) \neq g^T(x)] \leq 2me^{-T^2/2}$, so

$$\mathbb{P}[|g(x)| > t] \ge \frac{\|g^T\|^2 - t^2}{T^2 - t^2} - 2me^{-T^2/2}.$$

Thus by taking $T = \Theta(m)$ and $t = \Theta(1/m)$, we get that

$$\mathbb{P}[|g(x)| > \Theta(1/m)] \ge \Omega(1/m^4).$$

Thus finally we have

$$||f|| = \mathbb{E}[\tanh(g(x))^2] \ge \tanh^2(\Theta(1/m))\Omega(1/m^4) \ge \Omega(1/m^6),$$

since $tanh(x) \approx x - x^3$ for small x (by its Taylor series).

Lemma B.5 ([GGJ⁺20], Appendix A.1). For $ReLU^{T}(x) = min(T, ReLU(x))$,

$$\|\operatorname{ReLU} - \operatorname{ReLU}^T\|_{\mathcal{N}(0,1)} \le \sqrt{e^{-\frac{T^2}{2}} \left(T^2 + 1 - \frac{T}{\sqrt{2\pi}}\right)}.$$

Proof. Let $p(t) = \frac{1}{\sqrt{2\pi}}e^{-t^2/2}$ be the pdf of $\mathcal{N}(0,1)$. Then

$$\begin{aligned} \|\text{ReLU} - \text{ReLU}^T\|_{\mathcal{N}(0,1)}^2 &= \underset{t \sim \mathcal{N}(0,1)}{\mathbb{E}} \left[\left(\text{ReLU}(t) - \text{ReLU}^T(t) \right)^2 \right] \\ &= \int_T^\infty (t - T)^2 p(t) \ dt \\ &= \int_T^\infty t^2 p(t) \ dt - 2T \int_T^\infty t p(t) \ dt + T^2 \int_T^\infty p(t) \ dt \end{aligned}$$

Noting that p'(t) = -tp(t), we have

$$\int_{T}^{\infty} t^{2} p(t) dt = \int_{T}^{\infty} -t d(p(t))$$

$$= -t p(x) \Big|_{T}^{\infty} + \int_{T}^{\infty} p(t) dt \qquad \text{(integration by parts)}$$

$$= T p(T) + \underset{t \sim \mathcal{N}(0,1)}{\mathbb{P}} (t > T),$$

$$\int_{T}^{\infty} t p(t) dt = -p(t) \Big|_{T}^{\infty} = p(T),$$

$$\int_{T}^{\infty} p(t) dt = \underset{t \sim \mathcal{N}(0,1)}{\mathbb{P}} (t > T) \le e^{-\frac{T^{2}}{2}}.$$

The claim follows by algebra.

Lemma B.6. Consider f instantiated with $\psi = \tanh \ and \ \phi = \sigma$. Then $||f|| = e^{-O(\sqrt{m})}$.

Proof. Here the same approach as above becomes considerably simpler since $|g(x)| \leq 1$ always. The norm bound on g yields the following anti-concentration inequality:

$$\mathbb{P}[|g(x)| > t] \ge \frac{\|g\|^2 - t^2}{1 - t^2}.$$

In our case, taking $t=e^{-C\sqrt{m}}$ for sufficiently large C and using $\|g\|=e^{-O(\sqrt{m})}$ (Corollary B.3) yields

$$\mathbb{P}[|g(x)| > e^{-C\sqrt{m}}] = e^{-O(\sqrt{m})}.$$

Thus

$$||f|| = \mathbb{E}[\tanh(g(x))^2] \ge \tanh^2(e^{-C\sqrt{m}})e^{-O(\sqrt{m})} \ge e^{-O(\sqrt{m})}$$

since again $tanh(x) \approx x - x^3$ for small x.

C Approximate degree of ReLUs and sigmoids

Here we give estimates for the δ -approximate degree of ReLUs and sigmoids under the standard Gaussian using bounds on their Hermite coefficients. Recall that we consider units $\phi(w \cdot x)$ with $||w||_2 \le 1$. It is clear that for $\phi = \text{ReLU}$ and $\phi = \sigma$, the norm only increases monotonically with $||w||_2$, so for the purposes of analysis it suffices to consider exactly $||w||_2 = 1$.

It is not hard to show that whenever w is a unit vector, the total-degree-d Hermite weight of $\phi(w \cdot x)$ as $x \sim \mathcal{N}(0, I_n)$ is the same as that of the univariate $\phi(t)$ as $t \sim \mathcal{N}(0, 1)$. (A quick way of seeing this is to note that by rotational symmetry, we may assume WLOG that $w = e_1$, in which case the calculation is very straightforward.)

In what follows, we say $\widehat{\phi}_a$ are the Hermite coefficients of $\phi : \mathbb{R} \to \mathbb{R}$ if $\phi = \sum_a \widehat{\phi}_a H_a$, where the H_a are the normalized probabilists' Hermite polynomials. We use \widetilde{H}_a to denote the un-normalized (i.e. monic) Hermite polynomials. (Note that this is somewhat nonstandard notation.)

First we consider ReLUs.

Lemma C.1. $\widehat{\text{ReLU}}_0 = 1/\sqrt{2\pi}$, $\widehat{\text{ReLU}}_1 = 1/2$ and for $a \geq 2$, $\widehat{\text{ReLU}}_a = \frac{1}{\sqrt{2\pi a!}}(\widetilde{H}_a(0) + a\widetilde{H}_{a-2}(0))$. In particular, $\widehat{\text{ReLU}}_a = 0$ for odd $a \geq 3$ and $|\widehat{\text{ReLU}}_a| = \Theta(a^{-5/4})$ for even a.

Proof. We use the following standard recurrence relation: $\tilde{H}_{a+1}(x) = x\tilde{H}_a(x) - a\tilde{H}_{a-1}(x)$. For $a \geq 2$,

$$\widehat{\text{ReLU}}_{a} = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} \text{ReLU}(x) H_{a}(x) e^{-\frac{x^{2}}{2}} dx$$

$$= \frac{1}{\sqrt{2\pi a!}} \int_{0}^{\infty} x \tilde{H}_{a}(x) e^{-\frac{x^{2}}{2}} dx$$

$$= \frac{1}{\sqrt{2\pi a!}} \int_{0}^{\infty} (\tilde{H}_{a+1}(x) + a\tilde{H}_{a-1}(x)) e^{-\frac{x^{2}}{2}} dx$$

$$= \frac{1}{\sqrt{2\pi a!}} (\tilde{H}_{a}(0) + a\tilde{H}_{a-2}(0)).$$

Since $\tilde{H}_a(0) = 0$ for odd a, $\widehat{\text{ReLU}}_a = 0$ as well. For even a = 2b with $b \geq 2$, by standard expressions for $\tilde{H}_a(0)$, we have

$$\widehat{\text{ReLU}}_a = \frac{1}{\sqrt{2\pi(2b)!}} (\widetilde{H}_{2b}(0) + 2b\widetilde{H}_{2b-2}(0))
= \frac{1}{\sqrt{2\pi(2b)!}} \left((-1)^b \frac{(2b)!}{b!2^b} + 2b(-1)^{b-1} \frac{(2b-2)!}{(b-1)!2^{b-1}} \right)
= \frac{(-1)^b \sqrt{(2b)!}}{\sqrt{2\pi}b!2^b} \left(1 - \frac{2b}{2b-1} \right)
= \frac{(-1)^{b+1} \sqrt{(2b)!}}{\sqrt{2\pi}(2b-1)b!2^b}
= \frac{(-1)^{b+1}}{\sqrt{2\pi}(2b-1)(2b)^{1/4}}
= \frac{(-1)^{b+1}}{b^{5/4}}$$

Here the second inequality follows from the fact $\binom{n}{n/2} \approx \frac{2^{n/2}}{\sqrt{n}}$.

Corollary C.2. The δ -approximate degree of ReLU under $\mathcal{N}(0,1)$ is $O(\delta^{-4/3})$.

Proof. Let p denote the Hermite expansion of ReLU truncated at degree d. By the fact that $|\widehat{ReLU}_a| = \Theta(a^{-5/4})$ for even a (and 0 for odd a), we see that

$$||p - \text{ReLU}||^2 = \sum_{a>d} \widehat{\text{ReLU}}_a^2$$
$$= \sum_{\substack{a>d\\a \text{ even}}} \Theta(a^{-5/2})$$
$$= \Theta(d^{-3/2}).$$

For this to be at most δ^2 , we only need $d = O(\delta^{-4/3})$.

Now we turn to sigmoids. Let σ denote the standard sigmoid, i.e. the logistic function $\sigma(t) = 1/(1 + e^{-t})$.

Lemma C.3. For all sufficiently large a, $\hat{\sigma}_a = e^{-\Omega(\sqrt{a})}$.

Proof. Upper bounds on the Hermite coefficients of sigmoidal funtions are known to follow from classic results in the complex analysis of Hermite series [Hil40, Boy84]. We refer to [PSG19, Corollary F.7.1], where this computation is done for $\tanh'(x) = 1 - \tanh^2(x)$. The calculation is very similar for σ (in fact, σ is just an affine shift of \tanh).

Corollary C.4. The δ -approximate degree of σ under $\mathcal{N}(0,1)$ is $\tilde{O}(\log^2 1/\delta)$.

Proof. Let p denote the Hermite expansion of σ truncated at degree d. Observe that

$$\|\sigma - p\|^2 = \sum_{a>d} \widehat{\sigma}_a^2$$

$$= \sum_{a>d} e^{-\Omega(\sqrt{a})}$$

$$= \Theta(\sqrt{d}e^{-\Omega(\sqrt{d})}),$$

which is at most δ^2 for $d = \tilde{O}(\log^2 1/\delta)$.

D Frank-Wolfe convergence guarantee

Here we provide a self-contained proof of Theorem 2.5, restated here. In fact, we generalize the analysis to handle any constant factor approximation to the optimum, meaning that in the Frank–Wolfe subproblem of Algorithm 1, we only require

$$\langle s, -\nabla p(z_t) \rangle \ge \alpha \max_{s' \in \mathcal{Z}'} \langle s', -\nabla p(z_t) \rangle - \frac{1}{2} \delta \gamma_t C_p$$
 (10)

for some constant $\alpha \leq 1$. We closely follow [Jag13, Appendix A], noting the differences in our slightly more general setup (the standard setup has $\mathcal{Z}' = \mathcal{Z}$, and $\alpha = 1$).

Theorem D.1. Let $\mathcal{Z}' \subseteq \mathcal{Z}$ be convex sets, and let $p : \mathcal{Z} \to \mathbb{R}$ be a β -smoothly convex function. Let $C_p = \beta \operatorname{diam}(\mathcal{Z})^2$. Suppose that $z^* \in \mathcal{Z}'$ achieves $\min_{z' \in \mathcal{Z}'} p(z')$. For every t,

the iterates of Algorithm 1 (modified to work with Eq. (10)) satisfy

$$p(z_t) - p(z^*) \le \frac{2C_p}{\alpha^2(t+2)}(1+\delta).$$

Proof. Define the duality gap function $q: \mathcal{Z} \to \mathbb{R}$ as

$$q(z) = \max_{s \in \mathcal{Z}'} \langle z - s, \nabla p(z) \rangle.$$

Notice that q takes in any $z \in \mathcal{Z}$ but maximizes only over $s \in \mathcal{Z}'$. By convexity of p over \mathcal{Z} , we know that for all $z \in \mathcal{Z}, s \in \mathcal{Z}'$, $p(z) + \langle s - z, \nabla p(z) \rangle \leq p(s)$, meaning that $p(z) - p(s) \leq q(z)$. In particular, $p(z) - p(z^*) \leq q(z)$, so that q(z) always provides an upper bound on the gap between p(z) and $p(z^*)$ — this is weak duality.

Next we establish the following guarantee on the progress made in each step, which corresponds to Lemma 5 in Jaggi's proof.

Claim. Let the t^{th} step be $z_{t+1} = z_t + \gamma(s - z_t)$, where $z_t, z_{t+1}, s \in \mathcal{Z}$, $\gamma \in [0, 1]$ is arbitrary, and s satisfies

$$\langle s, -\nabla p(z_t) \rangle \ge \alpha \max_{s' \in \mathcal{Z}'} \langle s', -\nabla p(z_t) \rangle - \frac{1}{2} \delta \gamma C_p.$$

Then we have

$$p(z_{t+1}) \le p(z_t) - \alpha \gamma q(z_t) + \frac{\gamma^2}{2} C_p(1+\delta).$$

To see this, first note that because p is β -smoothly convex,

$$p(z_{t+1}) = p(z_t + \gamma(s - z_t))$$

$$\leq p(z_t) + \gamma \langle s - z_t, \nabla p(z_t) \rangle + \frac{\gamma^2}{2} C_p.$$

And from the way $s \in \mathcal{Z}$ was picked, we have

$$\langle s - z_t, -\nabla p(z_t) \rangle \ge \alpha \max_{s' \in \mathcal{Z}'} \langle s' - z_t, -\nabla p(z_t) \rangle - \frac{1}{2} \delta \gamma C_p$$
$$= \alpha q(z_t) - \frac{1}{2} \delta \gamma C_p.$$

The claim now follows.

As a consequence of the claim, we can say

$$p(z_{t+1}) - p(z^*) \le p(z_t) - p(z^*) - \gamma q(z_t) + \frac{\gamma^2}{2} C_p(1+\delta)$$

$$\le (1 - \alpha \gamma)(p(z_t) - p(z^*)) + \frac{\gamma^2}{2} C_p(1+\delta),$$

since $q(z_t) \ge p(z_t) - p(z^*)$ (weak duality). Taking $\gamma = \gamma_t = \frac{2}{\alpha(t+2)}$, the following bound can now by proven by induction on t:

$$p(z_t) - p(z^*) \le \frac{2}{\alpha^2(t+2)} C_p(1+\delta).$$

E Relationship between Boolean 0-1 loss and real-valued correlation loss

Let \mathcal{D} be a distribution on $\mathbb{R}^n \times \mathbb{R}$. Our lower bound applies against agnostic learners that satisfy Assumption 3.1, with a real-valued correlation guarantee, i.e. learners that learn a class \mathcal{H} by outputting $f: \mathbb{R}^n \to \mathbb{R}$ such that

$$\mathbb{E}_{(x,y)\sim\mathcal{D}}[f(x)y] \ge \max_{g\in\mathcal{H}} \mathbb{E}_{(x,y)\sim\mathcal{D}}[g(x)y] - \epsilon. \tag{11}$$

In the Boolean setting, where the labels are $\{\pm 1\}$ -valued, we have a distribution P on $\mathbb{R}^n \times \{\pm 1\}$. A learner is said to agnostically learn \mathcal{H} in terms of 0-1 loss if it is able to output $f: \mathbb{R}^n \to \{\pm 1\}$ such that

$$\mathbb{P}_{(a,b)\sim P}[f(a)\neq b] \leq \min_{g\in\mathcal{H}} \mathbb{P}_{(a,b)\sim P}[g(a)\neq b] + \epsilon,$$

or equivalently

$$\underset{(a,b)\sim P}{\mathbb{E}}[f(a)b] \ge \max_{g\in\mathcal{H}} \underset{(a,b)\sim P}{\mathbb{P}}[g(a)b] - \epsilon/2,$$

since $\mathbb{E}_{(a,b)\sim P}[f(a)b] = 1 - 2\mathbb{P}_{(a,b)\sim P}[f(a) \neq b]$. (The latter formulation has the benefit of making sense even for real-valued $f: \mathbb{R}^n \to \mathbb{R}$.)

It is not obvious that a learner L of the above kind (with a Boolean 0-1 loss guarantee) gives us a real-valued correlation loss guarantee, because it only knows how to operate on distributions P on $\mathbb{R}^n \times \{\pm 1\}$ (with Boolean labels), not distributions \mathcal{D} on $\mathbb{R}^n \times \mathbb{R}$ (with arbitrary real labels). Moreover, in the SQ setting, we must be able to translate L's queries to P, which are of the form $\phi: \mathbb{R}^n \times \{\pm 1\} \to \mathbb{R}$, into queries to \mathcal{D} . We claim that both of these difficulties can be gotten around. We will show that if \mathcal{D} has bounded labels, say in [-C, C], we can construct a distribution P on $\mathbb{R}^n \times \{\pm 1\}$ and simulate L on P to obtain a correlation loss guarantee wrt \mathcal{D} .

Indeed, let D denote the marginal of \mathcal{D} on \mathbb{R}^n ; for us, D is always $\mathcal{N}(0, I_n)$. Then P can be constructed simply as follows: draw $a \sim D$, and then randomly pick $b \in \{\pm 1\}$ such that $\mathbb{E}[b|a] = (\mathbb{E}_{(x,y)\sim\mathcal{D}}[y|x=a])/C$. (One could think of this as the "p-concept trick".) Equivalently, pick

$$b = \begin{cases} 1 & \text{with probability } \frac{1 + (\mathbb{E}_{(x,y) \sim \mathcal{D}}[y|x=a])/C}{2} \\ -1 & \text{otherwise} \end{cases}$$

One can easily see that for any $f: \mathbb{R}^n \to \mathbb{R}$,

$$\mathbb{E}_{(a,b)\sim P}[f(a)b] = \frac{1}{C} \mathbb{E}_{(x,y)\sim \mathcal{D}}[f(x)y],$$

so that using L to learn up to 0-1 error ϵ gives a correlation loss guarantee up to $C\epsilon/2$. It remains to show that we can indeed simulate L's queries to P using only SQ access to \mathcal{D} .

For any query $\phi: \mathbb{R}^n \times \{\pm 1\} \to \mathbb{R}$, observe that (since the marginal of P on \mathbb{R}^n is also D)

$$\begin{split} \underset{(a,b)\sim P}{\mathbb{E}}[\phi(a,b)] &= \underset{a\sim D}{\mathbb{E}}\left[\phi(a,1)\frac{1+(\mathbb{E}_{(x,y)\sim\mathcal{D}}[y|x=a])/C}{2} + \phi(a,-1)\frac{1-(\mathbb{E}_{(x,y)\sim\mathcal{D}}[y|x=a])/C}{2}\right] \\ &= \frac{1}{2}\underset{a\sim D}{\mathbb{E}}[\phi(a,1)+\phi(a,-1)] + \frac{1}{2C}\underset{(x,y)\sim\mathcal{D}}{\mathbb{E}}[(\phi(x,1)-\phi(x,-1)y]. \end{split}$$

This expression can be computed using two statistical queries to \mathcal{D} (or even just one, since we know the marginal D).

In our reduction (Theorem 4.1), we end up using the base learner on labeled distributions \mathcal{D} where the labels correspond to the loss functional's gradient; when using surrogate loss, the label for x is $\psi(f^*(x)) - \psi(f(x))$. We see that this is indeed bounded in [-2,2], since $\psi: \mathbb{R} \to [-1,1]$. Recall that in solving the Frank-Wolfe subproblem we needed to worry about simulating SQ access to this \mathcal{D} using only SQ access to the true $D_{\psi \circ f^*}$ (see Eq. (4) and surrounding discussion). Here we actually have a further layer: we need to simulate SQ access to P using SQ access to \mathcal{D} , itself simulated using actual SQ access to $D_{\psi \circ f^*}$. But it is easily verified that by the argument just outlined, no trouble arises here, and that one can in fact also "directly" simulate P using $D_{\psi \circ f^*}$ by the same argument as used for Eq. (4).

F Relationship between square loss and correlation loss for ReLUs

Let \mathcal{D} be a distribution on $\mathbb{R}^n \times \mathbb{R}$, and assume the labels are bounded in [-C, C]. Our lower bounds apply to agnostic learners that satisfy Assumption 3.1, with a guarantee in terms of correlation, where the output hypothesis f must satisfy

$$\underset{(x,y)\sim\mathcal{D}}{\mathbb{E}}[f(x)y] \ge \max_{q\in\mathcal{H}} \underset{(x,y)\sim\mathcal{D}}{\mathbb{E}}[g(x)y] - \epsilon.$$

But agnostic learning of real-valued functions is usually phrased in terms of square loss:

$$\mathbb{E}_{(x,y)\sim\mathcal{D}}[(f(x)-y)^2] \le \min_{g\in\mathcal{H}} \mathbb{E}_{(x,y)\sim\mathcal{D}}[(g(x)-y)^2] + \epsilon'.$$

Here we show that for the class of ReLUs, $\mathcal{H} = \mathcal{H}_{ReLU}$, an agnostic learner L with a square loss guarantee can be used to satisfy Assumption 3.1. Fundamentally, this amounts to working out a geometric relationship between distances and projections in our function space, and much of the following argument can be viewed as a somewhat careful elaboration of what, in the familiar Euclidean setup, is more easily visualized.

For simplicity, throughout this section we will scale the class \mathcal{H}_{ReLU} so that the maximum norm of any function is 1:

$$\mathcal{H} = \mathcal{H}_{\text{ReLU}} = \{ \pm \sqrt{2} \operatorname{ReLU}(u \cdot x) \mid ||u||_2 \le 1 \}.$$

An important property of this class is that we can always scale a function $h \in \mathcal{H}$ to have any desired norm in [0,1] without leaving the class. That is, for any nonzero $h \in \mathcal{H}$ and any $\lambda \in [0,1]$, $\frac{\lambda}{\|h\|}h \in \mathcal{H}$. This follows simply from the fact that $\|\text{ReLU}(u \cdot x)\| = \|u\|_2/\sqrt{2}$.

We can think of this as saying that \mathcal{H} is a norm-bounded section of a convex cone.

Let $f_{\mathsf{cmf}}(x) = \mathbb{E}[y|x]$. Let h_{sq} be a minimizer over all $h \in \mathcal{H}$ of the squared loss, $\mathbb{E}_{(x,y)\sim\mathcal{D}}[(h(x)-y)^2]$. An equivalent and more convenient view is that this is a minimizer of the squared distance $\|h-f_{\mathsf{cmf}}\|^2$, since

$$\|h - f_{\mathsf{cmf}}\|^2 = \|h\|^2 - 2\langle h, f_{\mathsf{cmf}} \rangle + \|f_{\mathsf{cmf}}\|^2 = \mathop{\mathbb{E}}_{\mathcal{D}}[(h(x) - y)^2] + \|f_{\mathsf{cmf}}\|^2 - \mathop{\mathbb{E}}_{\mathcal{D}}[y^2],$$

and the latter terms are independent of h. This view is particularly important since it, combined with the fact that \mathcal{H} is essentially a bounded convex cone, gives us an orthogonal projection theorem. Specifically, it is the case that the norm of h_{sq} must be the length of the projection of f_{cmf} onto the line λh_{sq} for $\lambda \in [0,1]$ (assuming this length is at most 1; otherwise, the norm is 1). In other words,

$$||h_{\mathsf{sq}}|| = \min\{\langle \frac{h_{\mathsf{sq}}}{||h_{\mathsf{sq}}||}, f_{\mathsf{cmf}}\rangle, 1\}.$$
(12)

This can be seen by asking: for what $\lambda \in [0,1]$ is $\|\frac{\lambda}{\|h_{sq}\|}h_{sq} - f_{cmf}\|$ minimized? (The point being that h_{sq} could be rescaled to have norm λ .) By writing this as

$$\|\frac{\lambda}{\|h_{\mathsf{sq}}\|}h_{\mathsf{sq}} - f_{\mathsf{cmf}}\|^2 = \left(\lambda - \langle \frac{h_{\mathsf{sq}}}{\|h_{\mathsf{sq}}\|}, f_{\mathsf{cmf}}\rangle\right)^2 + \|f_{\mathsf{cmf}}\|^2 - \langle \frac{h_{\mathsf{sq}}}{\|h_{\mathsf{sq}}\|}, f_{\mathsf{cmf}}\rangle^2,$$

the observation follows immediately.² This projection theorem also tells us that $h_{sq} = 0$ iff f_{cmf} has no projection onto any $h \in \mathcal{H}$, i.e. $\langle h, f_{cmf} \rangle = 0$ for all $h \in \mathcal{H}$.³

Let h_{cor} be a maximizer of the correlation, $\mathbb{E}_{(x,y)\sim\mathcal{D}}[h(x)y] = \langle h, f_{\mathsf{cmf}} \rangle$. We may clearly assume that h_{cor} has the maximum possible norm, which is 1. We claim that in fact, h_{cor} can be taken to be $h_{\mathsf{sq}}/\|h_{\mathsf{sq}}\|$ (assuming $h_{\mathsf{sq}} \neq 0$; otherwise, $h_{\mathsf{cor}} = 0$ as well since, as noted, this means $\langle h, f_{\mathsf{cmf}} \rangle = 0$ for all $h \in \mathcal{H}$). To see why, first assume $h_{\mathsf{sq}} \neq 0$ and use the fact that for any nonzero $h \in \mathcal{H}$, the square loss achieved by $\frac{\|h_{\mathsf{sq}}\|}{\|h\|}h$ (i.e. h scaled to have h_{sq} 's norm) cannot be better than that of h_{sq} itself. Thus by an algebraic manipulation we have

$$\begin{aligned} \|h_{\mathsf{sq}} - f_{\mathsf{cmf}}\|^2 &\leq \left\|\frac{\|h_{\mathsf{sq}}\|}{\|h\|}h - f_{\mathsf{cmf}}\right\|^2 \\ &\Longrightarrow \langle \frac{h_{\mathsf{sq}}}{\|h_{\mathsf{sq}}\|}, f_{\mathsf{cmf}}\rangle \geq \langle \frac{h}{\|h\|}, f_{\mathsf{cmf}}\rangle \geq \langle h, f_{\mathsf{cmf}}\rangle. \end{aligned}$$

Since this holds for any $h \in \mathcal{H}$, we may take $h_{cor} = h_{sq}/\|h_{sq}\|$.

Now suppose we have an agnostic learner in terms of square loss that returns h such that

$$\|h - f_{\mathsf{cmf}}\|^2 \leq \|h_{\mathsf{sq}} - f_{\mathsf{cmf}}\|^2 + \epsilon'.$$

For a suitable choice of ϵ' (depending on the final desired ϵ), we would like to say that $h/\|h\|$ achieves correlation that is ϵ -competitive with h_{cor} . Indeed, if $h_{\mathsf{sq}} = 0$ this is trivial, since as noted this means $\langle h, f_{\mathsf{cmf}} \rangle = 0$ for all $h \in \mathcal{H}$. Otherwise, by comparing $\frac{\|h\|}{\|h_{\mathsf{sq}}\|} h_{\mathsf{sq}}$ (i.e. h_{sq}

²Note that here we are assuming $\langle h_{sq}, f_{cmf} \rangle \geq 0$ WLOG, since otherwise we would consider $-h_{sq}$.

³For another way to see this, for any nonzero $h \in \mathcal{H}$, expand $\|\lambda h - f_{\mathsf{cmf}}\|^2 \ge \|0 - f_{\mathsf{cmf}}\|^2$ and let $\lambda \to 0$.

scaled to have h's norm) with h_{sq} itself, we may say that

$$||h - f_{\mathsf{cmf}}||^2 \le ||h_{\mathsf{sq}} - f_{\mathsf{cmf}}||^2 + \epsilon' \le \left\| \frac{||h||}{||h_{\mathsf{sq}}||} h_{\mathsf{sq}} - f_{\mathsf{cmf}} \right\|^2 + \epsilon'.$$

Some rearrangement gives

$$\langle \frac{h}{\|h\|}, f_{\mathsf{cmf}} \rangle \ge \langle \frac{h_{\mathsf{sq}}}{\|h_{\mathsf{sq}}\|}, f_{\mathsf{cmf}} \rangle - \frac{\epsilon'}{2\|h\|}
= \langle h_{\mathsf{cor}}, f_{\mathsf{cmf}} \rangle - \frac{\epsilon'}{2\|h\|},$$
(13)

showing that $h/\|h\|$ is $\frac{\epsilon'}{2\|h\|}$ -competitive with h_{cor} .

But an issue here is that ||h|| could be very small, or even zero. We claim that we can actually address this separately as an easy case: it implies that we are in a trivial situation in which even the 0 function performs fairly well, and so even the best possible correlation must be quite small.

Lemma F.1. Let h be such that $||h - f_{\mathsf{cmf}}||^2 \le ||h_{\mathsf{sq}} - f_{\mathsf{cmf}}||^2 + \epsilon'$. Suppose $||h|| \le \eta$. Then $\langle h_{\mathsf{cor}}, f_{\mathsf{cmf}} \rangle \le \sqrt{\epsilon' + 2C\eta}$. In particular, the 0 function is $\sqrt{\epsilon' + 2C\eta}$ -competitive with h_{cor} .

Proof. By Cauchy–Schwarz,

$$\|0 - f_{\mathsf{cmf}}\|^2 - \|h - f_{\mathsf{cmf}}\|^2 = 2\langle h, f_{\mathsf{cmf}} \rangle - \|f_{\mathsf{cmf}}\|^2 \le 2\|h\| \|f_{\mathsf{cmf}}\| \le 2C\eta,$$

where we use $||f_{cmf}|| \leq C$ since the labels are assumed to be bounded in [-C, C]. Thus

$$||0 - f_{\mathsf{cmf}}||^2 \le ||h - f_{\mathsf{cmf}}||^2 + 2C\eta \le ||h_{\mathsf{sg}} - f_{\mathsf{cmf}}||^2 + \epsilon' + 2C\eta.$$

On the other hand, by definition of h_{sq} ,

$$||h_{sq} - f_{cmf}||^2 \le ||0 - f_{cmf}||^2$$

Put together, this means that the 0 function achieves nearly the same square loss as h_{sq} :

$$||h_{\mathsf{sq}} - f_{\mathsf{cmf}}||^2 \le ||0 - f_{\mathsf{cmf}}||^2 \le ||h_{\mathsf{sq}} - f_{\mathsf{cmf}}||^2 + \epsilon' + 2C\eta. \tag{14}$$

This lets us conclude that $||h_{sq}||$ must be small:

$$||h_{sq}||^2 = ||f_{cmf}||^2 - ||h_{sq} - f_{cmf}||^2 + 2\langle h_{sq} - f_{cmf}, h_{sq} \rangle \le \epsilon' + 2C\eta,$$

where we use Eq. (14) and the fact that by can rewrite Eq. (12) as $||h_{sq}|| \leq \langle \frac{h_{sq}}{||h_{sq}||}, f_{cmf} \rangle$, or $\langle h_{sq} - f_{cmf}, h_{sq} \rangle \leq 0$. But now since $||h_{sq}|| \leq \sqrt{\epsilon' + 2C\eta} < 1$ (ϵ' and η will be picked sufficiently small), Eq. (12) boils down to saying that

$$\langle h_{\rm cor}, f_{\rm cmf} \rangle = \langle \frac{h_{\rm sq}}{\|h_{\rm sq}\|}, f_{\rm cmf} \rangle = \|h_{\rm sq}\| \leq \sqrt{\epsilon' + 2C\eta}.$$

We can now put everything together.

Theorem F.2. Suppose we have an agnostic learner L for \mathcal{H}_{ReLU} under \mathcal{D} with a square loss guarantee. Then L can be used to yield a correlation guarantee, i.e. to satisfy Assumption 3.1.

Proof. Run L with $\epsilon' = \Theta(\epsilon^3)$ to get h such that $\|h - f_{\mathsf{cmf}}\|^2 \leq \|h_{\mathsf{sq}} - f_{\mathsf{cmf}}\|^2 + \epsilon'$. By Lemma F.1, if $\|h\| \leq \eta = \Theta(\epsilon^2)$, then 0 is ϵ -competitive with h_{cor} . So we may assume that $\|h\| \geq \Theta(\epsilon^2)$. But then by Eq. (13), since now $\frac{\epsilon'}{2\|h\|} \leq \epsilon$, we get that $h/\|h\|$ is ϵ -competitive with h_{cor} .