

Copyright  
by  
Aravind Gollakota  
2023

The Dissertation Committee for Aravind Gollakota  
certifies that this is the approved version of the following dissertation:

**New Computational and Statistical Characterizations of  
Neural Network Learning**

**Committee:**

Adam R. Klivans, Supervisor

Pravesh K. Kothari

Qiang Liu

Eric Price

**New Computational and Statistical Characterizations of  
Neural Network Learning**

**by  
Aravind Gollakota**

**Dissertation**

Presented to the Faculty of the Graduate School of  
The University of Texas at Austin  
in Partial Fulfillment  
of the Requirements  
for the Degree of

**Doctor of Philosophy**

**The University of Texas at Austin  
August 2023**

## Acknowledgments

I owe a great deal to many wonderful people who made this PhD possible. Thanks go above all to my advisor, Adam Klivans, for welcoming me in and giving me the space to discover my own interests, for his easygoing approach and sense of humor in both research and life, for his excellent taste in problems, and for believing in me even through the lows and the lulls.

Collaboration is one of the great pleasures of research, and I feel truly fortunate and grateful to have been able to work with a stellar list of colleagues during my time in grad school: Sitan Chen, Yuval Dagan, Giannis Daras, Alex Dimakis, Surbhi Goel, Parikshit Gopalan, Zhihan Jin, Sushrut Karmalkar, Adam Klivans, Pravesh Kothari, Daniel Liang, Raghu Meka, Kulin Shah, Konstantinos Stavropoulos, and Arsen Vasilyan. Surbhi and Sushi showed me the ropes when I came in; and Kostas and Kulin came in not needing to be shown the ropes at all. I'm also very grateful to those who mentored me through my tentative first steps in research: Julia Chuzhoy, Karola Mészáros, István Miklós, David Steurer, and Madhur Tulsiani.

Austin has been a fantastic home for the last six years, and I will cherish the many memories from my time spent here with a fabulous group of friends, among them Akanksha, Akshay, Ananya, Daniel, Eric, Josh, Justin, Kevin, Kostas, Kulin, Matt, Meghna, Raghava, Ridwan, Sepideh, Shivam, Surbhi, and, of course, Sushi. To my friends from Cornell (Abhijit, Akshay, Athith, Jun Wei, Neha, Nivedita, Sandeep, Soham, Somrita, Sush, and many more), and from back home in India (Aditi, Athyuttam, Mihir, Prateek, Shubham, and many more): thank you guys for being there.

Heartfelt thanks go also to my extended family both near and far for their love and invaluable support, and especially to my uncles Rama and Ravi for helping to make the journey to these shores possible at all.

Finally, to my parents, for their unceasing love, encouragement, faith, and wisdom, and for everything else besides: thank you so much.

# Abstract

## New Computational and Statistical Characterizations of Neural Network Learning

Aravind Gollakota, PhD  
The University of Texas at Austin, 2023

SUPERVISOR: Adam R. Klivans

A foundational goal of machine learning theory is to characterize the inherent computational and statistical complexity of some of the most basic tasks in machine learning. In this thesis, we present new results concerning two such tasks in neural network learning and beyond.

First, we study the question of when efficient algorithms can achieve high test accuracy on labeled data known to be consistent with a simple neural network. We present a set of results establishing the surprising computational intractability of this problem even in the benign setting where the inputs are drawn from a Gaussian, and the labels are perfectly consistent with a simple two-hidden-layer or even one-hidden-layer neural network. These hardness results illuminate what types of problem assumptions are necessary for efficient algorithms for this problem to be possible at all.

Next, we investigate the problem of testing whether a learning algorithm has fit the data as well as its guarantee claims. This is a serious issue for agnostic supervised learning (i.e. supervised learning with no assumptions on the labels), where most efficient algorithms make simplifying distributional assumptions such as Gaussianity. But such assumptions can be hard to verify, meaning it can be hard to check whether

the learner has actually succeeded. The recent elegant model of testable learning addresses this issue by replacing such hard-to-verify distributional assumptions with efficiently testable ones. We present both a broad algorithmic framework as well as a full statistical characterization of this model.

# Table of Contents

|   |           |
|---|-----------|
| <b>Chapter 1: Introduction</b>  | <b>11</b> |
| 1.1 When can we achieve high test accuracy on data known to be consistent with a simple neural network? . . . . . | 11        |
| 1.2 What general characterizations of learning can we provide in the SQ model? . . . . .                          | 13        |
| 1.3 When can we be sure a learning algorithm has fit the data as well as possible? . . . . .                      | 14        |
| <b>Chapter 2: Superpolynomial Lower Bounds for Learning One-Layer Neural Networks using Gradient Descent</b>      | <b>17</b> |
| 2.1 Introduction . . . . .  | 17        |
| 2.1.1 Our Results . . . . .   | 17        |
| 2.1.2 Our techniques . . . . .  | 19        |
| 2.1.3 Related Work . . . . .  | 22        |
| 2.2 Preliminaries . . . . .   | 24        |
| 2.2.1 The SQ model . . . . .  | 25        |
| 2.3 Orthogonal Family of Neural Networks . . . . .  | 28        |
| 2.4 SQ Lower Bounds . . . . .   | 32        |
| 2.4.1 SQ Lower Bounds for Real-valued Functions . . . . .   | 32        |
| 2.4.2 SQ Lower Bounds for $\rho$ -concepts . . . . .  | 32        |
| 2.5 Experiments . . . . .   | 36        |
| <b>Chapter 3: Hardness of Noise-Free Learning for Two-Hidden-Layer Neural Networks</b>                            | <b>40</b> |
| 3.1 Introduction . . . . .  | 40        |
| 3.1.1 Discussion and Related Work . . . . .   | 45        |
| 3.1.2 Technical Overview . . . . .  | 47        |
| 3.2 Preliminaries . . . . .   | 50        |
| 3.2.1 Notation . . . . .  | 50        |
| 3.2.2 Learning models . . . . .   | 51        |
| 3.2.3 Learning with Rounding . . . . .  | 51        |
| 3.2.4 Partial assignments . . . . .   | 53        |
| 3.3 Compressing the Daniely–Vardi Lift . . . . .  | 54        |
| 3.3.1 The DV Lift . . . . .   | 56        |

|  |   |            |
|--|---|------------|
| 3.3.2  | Saving One Hidden Layer via Compressibility . . . . .                                   | 59         |
| 3.3.3  | Proofs of Lemmas 3.3.9 and 3.3.10 . . . . .   | 68         |
| 3.4  | Statistical Query Lower Bound . . . . .   | 73         |
| 3.4.1  | SQ lower bound via parities . . . . .   | 75         |
| 3.4.2  | SQ lower bound via the LWR functions . . . . .  | 76         |
| 3.5  | Cryptographic Hardness Based on LWR . . . . .   | 77         |
| 3.6  | Hardness of Learning using Label Queries . . . . .                                      | 79         |
| <b>Chapter 4: Statistical-Query Lower Bounds via Functional Gradients</b>                            |   | <b>81</b>  |
| 4.1  | Introduction . . . . .  | 81         |
| 4.1.1  | Our Approach . . . . .  | 83         |
| 4.1.2  | Related Work . . . . .  | 84         |
| 4.1.3  | Organization . . . . .  | 86         |
| 4.2  | Preliminaries . . . . .   | 86         |
| 4.2.1  | Statistical Query (SQ) Model . . . . .  | 87         |
| 4.2.2  | Convex Optimization Basics . . . . .  | 89         |
| 4.3  | Functional gradient descent . . . . .   | 90         |
| 4.3.1  | Frank–Wolfe using statistical queries . . . . .   | 93         |
| 4.4  | Functional gradient descent guarantees on surrogate loss . . . . .                      | 93         |
| 4.5  | Lower bounds on learning ReLUs, sigmoids, and halfspaces . . . . .                      | 97         |
| 4.6  | Lower bounds on learning general non-polynomial activations . . . . .                   | 99         |
| 4.7  | Lower bounds on learning monomials . . . . .  | 100        |
| 4.8  | Upper bounds on learning ReLUs and sigmoids . . . . .                                   | 102        |
| <b>Chapter 5: The Polynomial Method is Universal for Distribution-Free Correlational SQ Learning</b> |   | <b>105</b> |
| 5.1  | Introduction . . . . .  | 105        |
| 5.1.1  | Related and prior work . . . . .  | 109        |
| 5.2  | Preliminaries . . . . .   | 112        |
| 5.2.1  | Learning in the statistical query model . . . . .                                       | 113        |
| 5.2.2  | Correlational variance . . . . .  | 114        |
| 5.2.3  | Orthogonalizing distributions . . . . .   | 116        |
| 5.2.4  | Threshold and approximate degree . . . . .  | 117        |
| 5.2.5  | Pattern matrices . . . . .  | 118        |
| 5.3  | Main theorem: correlational variance bounds via orthogonalizing distributions . . . . . | 119        |



|   |  |            |
|---|--|------------|
| 5.4   | CSQ lower bounds for PAC learning in terms of threshold degree . . .                           | 121        |
| 5.5   | CSQ lower bounds for agnostic learning in terms of approximate degree                          | 123        |
| 5.5.1   | CSQ lower bounds on attribute-efficient agnostic learning of sparse halfspaces . . . . .       | 126        |
| 5.6   | Hardness of approximation in terms of approximate degree . . . . .                             | 127        |
| 5.7   | A hard class of $p$ -concepts under the uniform distribution . . . . .                         | 129        |
| <b>Chapter 6: A Moment-Matching Approach to Testable Learning and a New Characterization of Rademacher Complexity</b> |  | <b>134</b> |
| 6.1   | Introduction . . . . .   | 134        |
| 6.1.1   | Our results . . . . .  | 135        |
| 6.1.2   | Subsequent work and open questions . . . . .   | 145        |
| 6.1.3   | Other related work . . . . .   | 147        |
| 6.2   | Preliminaries . . . . .  | 149        |
| 6.2.1   | Notation and conventions . . . . .   | 149        |
| 6.2.2   | Learning models . . . . .  | 150        |
| 6.2.3   | Bounded independence and sandwiching polynomials over the hypercube . . . . .                  | 151        |
| 6.2.4   | Rademacher complexity . . . . .  | 152        |
| 6.3   | Duality . . . . .  | 154        |
| 6.4   | Testable learning via moment matching . . . . .  | 158        |
| 6.4.1   | Warm-up: testable learning over the hypercube via $k$ -wise independence . . . . .             | 158        |
| 6.4.2   | A general algorithm using moment matching . . . . .  | 159        |
| 6.5   | Testably learning functions of halfspaces over strictly subexponential distributions . . . . . | 161        |
| 6.5.1   | Moment closeness implies distribution closeness . . . . .                                      | 163        |
| 6.5.2   | Approximate low-degree moment matching fools functions of halfspaces . . . . .                 | 167        |
| 6.5.3   | Proof of Theorem 6.5.2 . . . . .   | 170        |
| 6.6   | Sample complexity of testable learning . . . . .   | 171        |
| 6.6.1   | Upper bound . . . . .  | 171        |
| 6.6.2   | Lower bound . . . . .  | 173        |
| 6.7   | Discussion . . . . .   | 178        |
| 6.7.1   | Implications for the uniform convergence paradigm . . . . .                                    | 178        |
| 6.7.2   | Implications for sandwiching degree . . . . .  | 181        |

|  |            |
|--|------------|
| <b>Appendix A: Superpolynomial Lower Bounds for Learning One-Layer Neural Networks using Gradient Descent</b>          | <b>183</b> |
| A.1 Bounding the function norms under the Gaussian   | 183        |
| A.1.1 ReLU Activation  | 185        |
| A.1.2 Sigmoid Activation   | 189        |
| A.1.3 General activations  | 192        |
| A.2 SQ lower bound for real-valued functions proof   | 192        |
| <b>Appendix B: Hardness of Noise-Free Learning for Two-Hidden-Layer Neural Networks</b>                                | <b>194</b> |
| B.1 Barriers for constructing $N_3$  | 194        |
| B.2 Supporting lemmas for Section 3.3  | 195        |
| B.3 SQ lower bound for the LWR functions   | 197        |
| <b>Appendix C: Statistical-Query Lower Bounds via Functional Gradients</b>   | <b>201</b> |
| C.1 SQ lower bound subtleties  | 201        |
| C.1.1 Relationships between parameters   | 201        |
| C.1.2 The dependence of the query lower bound on the error $\epsilon$ and the tolerance $\tau$                         | 202        |
| C.2 Bounding the function norms of the [DKKZ20] construction   | 203        |
| C.3 Approximate degree of ReLUs and sigmoids   | 210        |
| C.4 Frank–Wolfe convergence guarantee  | 213        |
| C.5 Relationship between Boolean 0-1 loss and real-valued correlation loss   | 214        |
| C.6 Relationship between square loss and correlation loss for ReLUs  | 216        |
| <b>Appendix D: A Moment-Matching Approach to Testable Learning and a New Characterization of Rademacher Complexity</b> | <b>221</b> |
| D.1 Proof of strong duality in Theorem 6.3.2   | 221        |
| <b>Bibliography</b>  | <b>225</b> |

# Chapter 1: Introduction

A foundational goal of machine learning theory is to characterize the complexity of some of the most basic tasks in machine learning. When can we achieve high test accuracy on data known to be consistent with a simple neural network or Boolean function? When can we be sure a learning algorithm has fit the data as well as possible? This thesis explores these and other related questions using tools from computational and statistical learning theory as well as complexity theory. The focus is on studying the inherent computational complexity of key learning problems. The goal is both to develop new algorithms and to show lower bounds that illuminate what types of problem assumptions are necessary for efficient algorithms to be possible at all.

## 1.1 When can we achieve high test accuracy on data known to be consistent with a simple neural network?

Consider for a moment the remarkable fact that given a large labeled training dataset of images, one can easily train a deep neural network to achieve extremely high accuracy at, say, detecting whether a fresh test image depicts a human face. What explains this performance? Certainly the fact that high test accuracy is achieved tells us that there *exists* a deep neural network nearly consistent with the labeled data; this might naturally lead us to study the approximation power of neural networks. Moreover, we have also managed to efficiently *find* such a network; this raises the important question of optimization for neural networks. This thesis adopts the perspective of *computational learning theory* [Val84, Vap98]: even in the ideal case where the data is *exactly* consistent with a simple neural network — i.e., given labeled data  $(x, y) \in \mathbb{R}^d \times \mathbb{R}$  drawn from a distribution  $D$  where  $y = f(x)$  for some unknown but simple neural network  $f$  — when is it possible to efficiently learn *any* hypothesis (not necessarily a neural network) that achieves high test accuracy? This is one of the

most basic definitions in learning theory, also known as the Probably Approximately Correct (PAC) model.

An especially useful model for such problems is the Statistical Query (SQ) model [Kea98, BFJ<sup>+</sup>94, Rey20]. SQ learners are abstractions of many common learning algorithms (including gradient descent) that access the data exclusively through statistical queries of the form  $\mathbb{E}_{(x,y) \sim D}[\phi(x,y)] \pm \tau$  for some bounded query function  $\phi$  and tolerance  $\tau$ . In recent years, there has been a substantial body of work on showing lower bounds and statistical-computational tradeoffs for many supervised and unsupervised learning problems (see e.g. [DKS17, BB20, BBH<sup>+</sup>20] and references therein), and for such problems SQ lower bounds are viewed as strong evidence of the limits of computationally efficient learners.

In Chapter 2, we show that this problem is hard for a broad subclass of SQ learners known as correlational SQ (CSQ) learners (which include gradient descent on squared loss) even in the benign setting where the distribution over inputs is the standard Gaussian and  $f$  is just a one-hidden-layer network of the form  $f(x) = \sum_{i=1}^m a_i \sigma(w_i \cdot x)$ . This result, however, leaves open the question of ruling out *all* SQ learners (as opposed to just the CSQ subclass), and prevailing wisdom (see e.g. [VW19a]) actually suggested that this might not even be possible for a real-valued function family in the absence of noise. In Chapter 3, we answer this question and show the first *general* SQ lower bounds for learning noise-free *two*-hidden-layer ReLU networks over the Gaussian. In fact, we show that the same hard family of two-hidden-layer networks is also hard for *any* efficient learner (not necessarily SQ) under a cryptographic assumption. We also show the first hardness result for learning constant-depth networks using *label queries*.

On the other end of the spectrum from the noise-free setting is the agnostic setting, where the input distribution is still Gaussian but the labels may be arbitrarily noisy. For this setting, we show in Chapter 4 that the problem becomes harder still: merely finding a hypothesis competitive with the best-fitting *single ReLU* is beyond

any efficient SQ algorithm. Our proof introduces a novel algorithm for lifting a learner for single neurons to a learner for one-hidden-layer networks using Frank–Wolfe functional gradient descent.

Taken together, these constitute the strongest lower bounds to date on learning shallow neural networks over the Gaussian. From an algorithm designer’s standpoint, the value of these lower bounds is in showing what types of problem assumptions it is necessary to make for efficient algorithms to be possible at all.

## 1.2 What general characterizations of learning can we provide in the SQ model?

Moving beyond neural network learning, it is natural to wonder whether one can provide characterizations of the learnability of general Boolean concept classes, at least in the SQ model. In the distribution-free setting (in which the input distribution is unknown and arbitrary), the only known efficient learners for almost all rich classes are based on the so-called “polynomial method” [HS07]. This method relies on approximating Boolean functions by low-degree polynomials (in various senses depending on the precise setting), and then using simple learners for these polynomial representations. Is this the best we can do?

Remarkably, there is evidence suggesting that the answer is yes: there exist SQ lower bounds that essentially match the upper bounds given by the polynomial method. These lower bounds follow as consequences of various powerful results due to [Fel08, Fel12, She08b, She11], leveraging connections to other fields such as communication complexity, and hold primarily for the subclass of correlational SQ (CSQ) learners. While most known PAC algorithms fall into the SQ framework, it is worth pointing out that CSQ learners are only a subset of SQ learners, and there do exist results showing separations between CSQ and SQ learnability [Fel08]. Nevertheless, such CSQ lower bounds do constitute important evidence of hardness of learning in the general PAC model.

For such a basic result, alternative proofs are always useful, and in Chapter 5 we provide a new proof of the optimality of the polynomial method in the CSQ model. Our formal lower bound for noise-free/realizable learning is stated in terms of the so-called threshold degree, namely the minimum degree required for polynomials to approximate a given concept class in 0/1 loss (i.e. in sign), and for agnostic learning in terms of the approximate degree, namely the minimum degree required for polynomials to approximate a class pointwise. To our knowledge, the precise statements we give had not appeared in this form before, and help to further clarify the role of the polynomial method in computational learning theory. Moreover, our proofs are simple and largely self-contained.

### 1.3 When can we be sure a learning algorithm has fit the data as well as possible?

For many challenging problems in machine learning theory, simplifying assumptions are made in search of efficient algorithms. But an underappreciated fact about such assumptions is that if they cannot be efficiently verified, then for some problems it can be hard to verify *whether the learning algorithm has actually succeeded*. This is a significant issue for the usability as well as interpretability of these algorithms.

An important example is agnostic supervised learning, where the learner’s goal is produce a hypothesis with error close to the optimal error  $\text{opt}(\mathcal{C}, D)$  achievable by any function in a function class  $\mathcal{C}$  over a labeled distribution  $D$ . To make this problem tractable, it is common to make distributional assumptions such as Gaussianity. But these assumptions can be hard or impossible to verify given only samples from  $D$ , and because we have no prior knowledge of  $\text{opt}(\mathcal{C}, D)$ , it is not clear how to verify the learner’s error guarantee.

The recent elegant model of *testable learning* [RV23] addresses this issue by requiring that the learner succeed whenever an associated efficient tester accepts the

unknown distribution. Moreover, the tester is required to accept whenever the original distributional assumption that is being replaced is indeed satisfied.

In Chapter 6, we establish some foundational results for this model. We show a powerful algorithmic framework that yields testable learners matching the bounds of the optimal ordinary agnostic learners for many important classes. Our framework works for any function class that admits *sandwiching polynomials* with bounded coefficients, and may be viewed as an analog of the universal approach to agnostic learning based on  $L^1$  polynomial regression of [KKMS08]. The associated tester is simple and relies entirely on checking low-degree empirical moments. On the sample complexity front, we show a nearly-tight characterization of the sample complexity of testable learning in terms of Rademacher complexity, a well-studied quantity in statistical learning theory [BBL03]. Intriguingly, this means uniform convergence is both necessary and sufficient for testable learning, yielding a fundamental separation from ordinary distribution-specific agnostic learning. This is the first natural model-based characterization of Rademacher complexity. This separation also arguably gives us a new lens on understanding generalization in modern overparameterized models [BMR21, Bel21].

## Bibliographic note

This thesis is based primarily on the following published papers, appearing with the permission of the coauthors:

- Chapter 2 is based on joint work with Surbhi Goel, Zhihan Jin, Sushrut Karmalkar, and Adam Klivans. This work appeared in ICML 2020 [GGJ<sup>+</sup>20].
- Chapter 3 is based on joint work with Sitan Chen, Raghu Meka, and Adam Klivans. This work appeared in NeurIPS 2022 [CGMK22].
- Chapter 4 is based on joint work with Surbhi Goel and Adam Klivans. This work appeared in NeurIPS 2020 [GGK20].

- Chapter 6 is based on joint work with Pravesh Kothari and Adam Klivans. This work appeared in STOC 2023 [[GKK23](#)].



# Chapter 2: Superpolynomial Lower Bounds for Learning One-Layer Neural Networks using Gradient Descent

## 2.1 Introduction

A major challenge in the theory of deep learning is to understand when gradient descent can efficiently learn simple families of neural networks. The associated optimization problem is nonconvex and well known to be computationally intractable in the worst case. For example, cyphertexts from public-key cryptosystems can be encoded into a training set labeled by simple neural networks [KS09], implying that the corresponding learning problem is as hard as breaking cryptographic primitives. These hardness results, however, rely on discrete representations and produce relatively unrealistic joint distributions.

### 2.1.1 Our Results

In this chapter we give the first superpolynomial lower bounds for learning neural networks using gradient descent in arguably the simplest possible setting: we assume the marginal distribution is a spherical Gaussian, the labels are noiseless and are exactly equal to the output of a one-layer neural network (a linear combination of say ReLU or sigmoid activations), and the goal is to output a classifier whose test error (measured by square-loss) is small. We prove—unconditionally—that gradient descent fails to produce a classifier with small square-loss if it is required to run in polynomial time in the dimension. Our lower bound depends only on the algorithm used (gradient descent) and not on the architecture of the underlying classifier. That is, our results imply that current popular heuristics such as running gradient descent on an overparameterized network (for example, working in the NTK regime [JHG18]) will require superpolynomial time to achieve small test error.

**Statistical Queries.** We prove our lower bounds in the now well-studied statistical query (SQ) model of [Kea98] that captures most learning algorithms used in practice. For a loss function  $\ell$  and a hypothesis  $h_\theta$  parameterized by  $\theta$ , the true population loss with respect to joint distribution  $D$  on  $X \times Y$  is given by  $\mathbb{E}_{(x,y) \sim D}[\ell(h_\theta(x), y)]$ , and the gradient with respect to  $\theta$  is given by  $\mathbb{E}_{(x,y) \sim D}[\ell^\theta(h_\theta(x), y) \nabla_\theta h_\theta(x)]$ . In the SQ model, we specify a query function  $\phi(x, y)$  and receive an estimate of  $\int \mathbb{E}_{(x,y) \sim D}[\phi(x, y)]$  to within some tolerance parameter  $\tau$ . An important special class of queries are correlational or inner-product queries, where the query function  $\phi$  is defined only on  $X$  and we receive an estimate of  $\int \mathbb{E}_{(x,y) \sim D}[\phi(x) \cdot y]$  within some tolerance  $\tau$ . It is not difficult to see that (1) the gradient of a population loss can be approximated to within  $\tau$  using statistical queries of tolerance  $\tau$  and (2) for square-loss only inner-product queries are required.

Since the convergence analysis of gradient descent holds given sufficiently strong approximations of the gradient, lower bounds for learning in the SQ model [Kea98, BFJ<sup>+</sup>94, Szö09, Fel12, Fel17] directly imply unconditional lower bounds on the running time for gradient descent to achieve small error. We give the first super-polynomial lower bounds for learning one-layer networks with respect to any Gaussian distribution for any SQ algorithm that uses inner product queries:

**Theorem 2.1.1** (informal). *Let  $C$  be a class of real-valued concepts defined by one-layer single-output neural networks with input dimension  $n$  and  $m$  hidden units (ReLU or sigmoid); i.e., functions of the form  $f(x) = \sum_{i=1}^m a_i \sigma(w_i \cdot x)$ . Then learning  $C$  under the standard Gaussian  $N(0, I_n)$  in the SQ model with inner-product queries requires  $n^{\Omega(\log m)}$  queries for any tolerance  $\tau = n^{-\Omega(\log m)}$ .*

In particular, this rules out any approach for learning one-layer neural networks in polynomial-time that performs gradient descent on any polynomial-size classifier with respect to square-loss or logistic loss. For classification, we obtain significantly stronger results and rule out general SQ algorithms that run in polynomial-time (e.g.,

gradient descent with respect to any polynomial-size classifier and any polynomial-time computable loss). In this setting, our labels are  $f \in \{-1, 1\}^n$  and correspond to the *softmax* of an unknown one-layer neural network. We prove the following:

**Theorem 2.1.2** (informal). *Let  $\mathcal{C}$  be a class of real-valued concepts defined by a one-layer neural network in  $n$  dimensions with  $m$  hidden units (ReLU or sigmoid) feeding into any odd, real-valued output node with range  $[-1, 1]$ . Let  $D^\theta$  be a distribution on  $\mathbb{R}^n \times \{-1, 1\}^n$  such that the marginal on  $\mathbb{R}^n$  is the standard Gaussian  $N(0, I_n)$ , and  $E[Y|X] = c(X)$  for some  $c \in \mathcal{C}$ . For some  $b, C > 0$  and  $\epsilon = Cm^{-b}$ , outputting a classifier  $f : \mathbb{R}^n \rightarrow \{-1, 1\}^n$  with  $P_{(X,Y)}[f(X) \neq Y] \leq 1/2 + \epsilon$  requires  $n = \Omega(m^{1/b} \log m)$  statistical queries of tolerance  $\epsilon$ .*

The above lower bound for classification rules out the commonly used approach of training a polynomial-size, real-valued neural network using gradient descent (with respect to any polynomial-time computable loss) and then taking the sign of the output of the resulting network.

### 2.1.2 Our techniques

At the core of all SQ lower bounds is the construction of a family of functions that are pairwise approximately orthogonal with respect to the underlying marginal distribution. Typically, these constructions embed  $2^n$  parity functions over the discrete hypercube  $\{-1, 1\}^n$ . Since parity functions are perfectly orthogonal, the resulting lower bound can be quite strong. Here we wish to give lower bounds for more natural families of distributions, namely Gaussians, and it is unclear how to embed parity.

Instead, we use an alternate construction. For activation functions  $\phi, \psi : \mathbb{R} \rightarrow \mathbb{R}$ , define

$$f_S(x) = \psi \left( \sum_{w \in \{-1, 1\}^k} \chi(w) \phi \left( \frac{w \cdot x_S}{k} \right) \right).$$

Enumerating over every  $S \subseteq [n]$  of size  $k$  gives a family of functions of size  $n^{O(k)}$ . Here  $x_S$  denotes the vector of  $x_i$  for  $i \in S$  (typically we choose  $k = \log m$  to produce a family of one-layer neural networks with  $m$  hidden units). Each of the  $2^k = m$  inner weight vectors are all of unit norm, and all of the  $m$  outer weights have absolute value one. Note also that our construction uses activations with zero bias term.

We give a complete characterization of the class of nonlinear activations for which these functions are orthogonal. In particular, the family is orthogonal for any activation with a nonzero Hermite coefficient of degree  $k$  or higher.

Apart from showing orthogonality, we must also prove that functions in these classes are nontrivial (i.e., are not exponentially close to the constant zero function). This reduces to proving certain lower bounds on the norms of one-layer neural networks. The analysis requires tools from Hermite and complex analysis.

**SQ Lower Bounds for Real-Valued Functions.** Another major challenge is that our function family is real-valued as opposed to boolean. Given an orthogonal family of (deterministic) *boolean* functions, it is straightforward to apply known results and obtain general SQ lower bounds for learning with respect to 0/1 loss. For the case of real-valued functions, the situation is considerably more complicated. For example, the class of orthogonal Hermite polynomials on  $n$  variables of degree  $d$  has size  $n^{O(d)}$ , yet there *is* an SQ algorithm due to [APVZ14] that learns this class with respect to the Gaussian distribution in time  $2^{O(d)}$ . More recent work due to [ADHV19] shows that Hermite polynomials can be learned by an SQ algorithm in time polynomial in  $n$  and  $\log d$ .

As such, it is *impossible* to rule out general polynomial-time SQ algorithms for learning real-valued functions based solely on orthogonal function families. Fortunately, it is not difficult to see that the SQ reductions due to [Szö09] hold in the real-valued setting as long as the learning algorithm uses only *inner-product queries* (and the norms of the functions are sufficiently large). Since performing gradient

descent with respect to square-loss or logistic loss can be implemented using inner-product queries, we obtain our first set of desired results<sup>1</sup>.

Still, we would like rule out *general* SQ algorithms for learning simple classes of neural networks. To that end, we consider the *classification* problem for one-layer neural networks and output labels after performing a *softmax* on a one-layer network. Concretely, consider a distribution on  $\mathbb{R}^n \times \{-1, 1\}$  where  $\mathbb{E}[Y|X] = \sigma(c(X))$  for some  $c \in \mathcal{C}$  and  $\sigma : \mathbb{R} \rightarrow [-1, 1]$  (for example,  $\sigma$  could be  $\tanh$ ). We describe two goals. The first is to estimate the conditional mean function, i.e., output a classifier  $h$  such that  $\mathbb{E}[(h(x) - c(x))^2] \leq \epsilon$ . The second is to directly minimize classification loss, i.e., output a boolean classifier  $h$  such that  $\mathbb{P}_{X,Y} [h(X) \neq Y] \leq 1/2 + \epsilon$ .

We give superpolynomial lower bounds for both of these problems in the general SQ model by making a new connection to *probabilistic concepts*, a learning model due to [KS94a]. Our key theorem gives a superpolynomial SQ lower bound for the problem of *distinguishing* probabilistic concepts induced by our one-layer neural networks from truly random labels. A final complication we overcome is that we must prove orthogonality and norm bounds on one-layer neural networks that have been composed with a nonlinear activation (e.g.,  $\tanh$ ).

**SGD and Gradient Descent Plus Noise.** It is easy to see that our results also imply lower bounds for algorithms where the learner adds noise to the estimate of the gradient (e.g., Langevin dynamics). On the other hand, for technical reasons, it is known that SGD is *not* a statistical query algorithm (because it examines training points individually) and does not fall into our framework. That said, recent work by [AS20] shows that SGD is *universal* in the sense that it can encode all polynomial-time learners. This implies that proving unconditional lower bounds for SGD would give a proof that  $\text{P} \neq \text{NP}$ . Thus, we cannot hope to prove unconditional lower bounds on SGD (unless we can prove  $\text{P} \neq \text{NP}$ ).

---

<sup>1</sup>The algorithms of [APVZ14] and [ADHV19] do not use inner-product queries.

### 2.1.3 Related Work

There is a large literature of results proving hardness results (or unconditional lower bounds in some cases) for learning various classes of neural networks [BR89, Vu98, KS09, LSSS14, GKKT17b].

The most relevant prior work is due to [SVWX17], who addressed learning one-layer neural networks under logconcave distributions using Lipschitz queries. Specifically, let  $n$  be the input dimension, and let  $m$  be the number of hidden  $s$ -Lipschitz sigmoid units. For  $m = \tilde{O}(s^{\rho} \bar{n})$ , they construct a family of neural networks such that any learner using  $\lambda$ -Lipschitz queries with tolerance greater than  $\Omega(1/(s^2 n))$  needs at least  $2^{\binom{n}{2}}/(\lambda s^2)$  queries.

Roughly speaking, their lower bounds hold for  $\lambda$ -Lipschitz queries due to the composition of their one-layer neural networks with a  $\delta$ -function in order make the family more “boolean.” Because of their restriction on the tolerance parameter, they cannot rule out gradient descent with large batch sizes. Further, the slope of the activations they require in their constructions scales inversely with the Lipschitz and tolerance parameters.

To contrast with [SVWX17], note that our lower bounds hold for any inverse-polynomial tolerance parameter (i.e., *will* hold for polynomially-large batch sizes), do not require a Lipschitz constraint on the queries, and use only standard 1-Lipschitz ReLU and/or sigmoid activations (with zero bias) for the construction of the hard family. Our lower bounds are typically quasipolynomial in the number of hidden units; improving this to an exponential lower bound is an interesting open question. Both of our models capture square-loss and logistic loss.

In terms of techniques, [SVWX17] build an orthogonal function family using univariate, periodic “wave” functions. Our construction takes a different approach, adding and subtracting activation functions with respect to overlapping “masks.” Finally, aside from the (black-box) use of a theorem from complex analysis, our construction and analysis are considerably simpler than the proof in [SVWX17].

A follow-up work [VW19b] gave SQ lower bounds for learning classes of degree  $d$  orthogonal polynomials in  $n$  variables with respect to the uniform distribution on the unit sphere (as opposed to Gaussians) using inner product queries of bounded tolerance (roughly  $1/n^d$ ). To obtain superpolynomial lower bounds, each function in the family requires superpolynomial description length (their polynomials also take on very small values,  $1/n^d$ , with high probability).

Shamir [Sha18a] (see also the related work of [SSSS17]) proves hardness results (and lower bounds) for learning neural networks using gradient descent with respect to square-loss. His results are separated into two categories: (1) hardness for learning “natural” target families (one layer ReLU networks) or (2) lower bounds for “natural” input distributions (Gaussians). We achieve lower bounds for learning problems with *both* natural target families and natural input distributions. Additionally, our lower bounds hold for any nonlinear activations (as opposed to just ReLUs) and for broader classes of algorithms (SQ).

Recent work due to [GKK19] gives hardness results for learning a ReLU with respect to Gaussian distributions. Their results require the learner to output a single ReLU as its output hypothesis and require the learner to succeed in the agnostic model of learning. [KK14b] prove hardness results for learning a threshold function with respect to Gaussian distributions, but they also require the learner to succeed in the agnostic model. Very recent work due to Daniely and Vardi [DV20a] gives hardness results for learning randomly chosen two-layer networks. The hard distributions in their case are not Gaussians, and they require a nonlinear clipping output activation.

**Positive Results.** Many recent works give algorithms for learning one-layer ReLU networks using gradient descent with respect to Gaussians under various assumptions [ZSJ<sup>+</sup>17a, ZPS17, BG17a, ZYWG19a] or use tensor methods [JSA15, GLM18a].

These results depend on the hidden weight vectors being sufficiently orthogonal, or the coefficients in the second layer being positive, or both. Our lower bounds explain why these types of assumptions are necessary.

**Independent Work.** Independently, Diakonikolas et al. [DKKZ20] have given stronger correlational SQ lower bounds for the same class of functions with respect to the Gaussian distribution. Their bounds are exponential in the number of hidden units while ours is quasipolynomial. We can plug in their result and obtain exponential general SQ lower bounds for the associated probabilistic concept using our framework.

## 2.2 Preliminaries

We use  $[n]$  to denote the set  $\{1, \dots, n\}$ , and  $S \subseteq_k T$  to indicate that  $S$  is a  $k$ -element subset of  $T$ . We denote euclidean inner products between vectors  $u$  and  $v$  by  $\langle u, v \rangle$ . We denote the element-wise product of vectors  $u$  and  $v$  by  $u \odot v$ , that is,  $u \odot v$  is the vector  $(u_1 v_1, \dots, u_n v_n)$ .

Let  $X$  be an arbitrary domain, and let  $D$  be a distribution on  $X$ . Given two functions  $f, g : X \rightarrow \mathbb{R}$ , we define their  $L_2$  inner product with respect to  $D$  to be  $\langle f, g \rangle_D = \mathbb{E}_D[fg]$ . The corresponding  $L_2$  norm is given by  $\|f\|_D = \sqrt{\langle f, f \rangle_D} = \sqrt{\mathbb{E}_D[f^2]}$ .

A real-valued concept on  $\mathbb{R}^n$  is a function  $c : \mathbb{R}^n \rightarrow \mathbb{R}$ . We denote the induced labeled distribution on  $\mathbb{R}^n \times \mathbb{R}$ , i.e. the distribution of  $(x, c(x))$  for  $x \sim D$ , by  $D_c$ . A probabilistic concept, or  $p$ -concept, on  $X$  is a concept that maps each point  $x$  to a random  $\{-1, 1\}$ -valued label in such a way that  $\mathbb{E}[Y|X] = c(X)$  for a fixed function  $c : \mathbb{R}^n \rightarrow [-1, 1]$ , known as the conditional mean function. Given a distribution  $D$  on the domain, we abuse  $D_c$  to denote the induced labeled distribution on  $X \times \{-1, 1\}$  such that the marginal distribution on  $\mathbb{R}^n$  is  $D$  and  $\mathbb{E}[Y|X] = c(X)$  (equivalently the label is  $+1$  with probability  $\frac{1+c(x)}{2}$  and  $-1$  otherwise).



### 2.2.1 The SQ model

A statistical query is specified by a query function  $h : \mathbb{R}^d \rightarrow \mathbb{R}$ . The SQ model allows access to an SQ oracle that accepts a query  $h$  of specified tolerance  $\tau$ , and responds with a value in  $[\mathbb{E}_{x,y}[h(x,y)] - \tau, \mathbb{E}_{x,y}[h(x,y)] + \tau]$ .<sup>2</sup> To disallow arbitrary scaling, we will require that for each  $y$ , the function  $x \mapsto h(x,y)$  has norm at most 1. In the real-valued setting, a query  $h$  is called a correlational or inner product query if it is of the form  $h(x,y) = g(x) \cdot y$  for some function  $g$ , so that  $\mathbb{E}_{D_c}[h] = \mathbb{E}_D[gc] = \langle g, c \rangle_D$ . Here we assume  $\|g\| \leq 1$  when stating lower bounds, again to disallow arbitrary scaling.

Gradient descent with respect to squared loss is captured by inner product queries, since the gradient is given by

$$\begin{aligned} \mathbb{E}_{x,y} [r_\theta(h_\theta(x) - y)^2] &= \mathbb{E}_{x,y} [2(h_\theta(x) - y)r_\theta h_\theta(x)] \\ &= 2 \mathbb{E}_x [h_\theta(x)r_\theta h_\theta(x)] \\ &= 2 \mathbb{E}_{x,y} [y r_\theta h_\theta(x)]. \end{aligned}$$

Here the first term can be estimated directly using knowledge of the distribution, while the latter is a vector each of whose elements is an inner product query.

We now formally define the learning problems we consider.

**Definition 2.2.1** (SQ learning of real-valued concepts using inner product queries). Let  $\mathcal{C}$  be a class of  $p$ -concepts over a domain  $X$ , and let  $D$  be a distribution on  $X$ . We say that a learner learns  $\mathcal{C}$  with respect to  $D$  up to  $L_2$  error  $\epsilon$  using inner product queries (equivalently squared loss  $\epsilon^2$ ) if, given only SQ oracle access to  $D_c$  for some unknown  $c \in \mathcal{C}$ , and using only inner product queries, it is able to output  $\tilde{c} : X \rightarrow [-1, 1]$  such that  $\|c - \tilde{c}\|_D \leq \epsilon$ .

---

<sup>2</sup>In the SQ literature, this is referred to as the STAT oracle. A variant called VSTAT is also sometimes used, known to be equivalent up to small polynomial factors [Fel17]. While it makes no substantive difference to our superpolynomial lower bounds, our arguments can be extended to VSTAT as well.

For the classification setting, we consider two different notions of learning  $p$ -concepts. One is learning the target up to small  $L_2$  error, to be thought of as a strong form of learning. The other, weaker form, is achieving a nontrivial inner product (i.e. unnormalized correlation) with the target. We prove lower bounds on both in order to capture different learning goals.

**Definition 2.2.2** (SQ learning of  $p$ -concepts). Let  $\mathcal{C}$  be a class of  $p$ -concepts over a domain  $X$ , and let  $D$  be a distribution on  $X$ . We say that a learner learns  $\mathcal{C}$  with respect to  $D$  up to  $L_2$  error  $\epsilon$  if, given only SQ oracle access to  $D_c$  for some unknown  $c \in \mathcal{C}$ , and using arbitrary queries, it is able to output  $\tilde{c} : X \rightarrow [-1, 1]$  such that  $\|c - \tilde{c}\|_D \leq \epsilon$ . We say that a learner weakly learns  $\mathcal{C}$  with respect to  $D$  with advantage  $\epsilon$  if it is able to output  $\tilde{c} : X \rightarrow [-1, 1]$  such that  $\langle \tilde{c}, c \rangle_D \geq \epsilon$ .

Note that the best achievable advantage is  $\mathbb{E}_x \mathbb{E}_D [|\langle c, f \rangle|]$ , achieved by  $\tilde{c}(x) = \text{sign}(c(x))$ . Note also that  $\|c\|_D^2 = \mathbb{E}_D [c^2] = \|c\|_D$ , and therefore a norm lower bound on functions in  $\mathcal{C}$  implies an upper bound on the achievable advantage.

*Remark 2.2.3* (Learning with  $L_2$  error implies weak learning). If the functions in our class satisfy a norm lower bound, say  $\|c\|_D^2 \geq (1 + \alpha)\epsilon^2$ , then a simple calculation shows that learning with  $L_2$  error  $\epsilon$  implies weak learning with advantage  $\alpha\epsilon^2/2$ .

Our definition of weak learning also captures the standard boolean sense of weak learning, in which the learner is required to output a boolean hypothesis with 0/1 loss bounded away from 1/2. Indeed, by an easy calculation, the 0/1 loss of a function  $f : X \rightarrow [-1, 1]$  satisfies

$$\mathbb{P}_{(x,y) \sim D_c} [f(x) \neq y] = \frac{1}{2} - \frac{\langle f, c \rangle_D}{2}.$$

The difficulty of learning a concept class in the SQ model is captured by a parameter known as the statistical dimension of the class.

**Definition 2.2.4** (Statistical dimension). Let  $\mathcal{C}$  be a concept class of either real-valued concepts or  $p$ -concepts (i.e. their corresponding conditional mean functions)

on a domain  $X$ , and let  $D$  be a distribution on  $X$ . The (un-normalized) *correlation* of two concepts  $c, c' \in \mathcal{C}$  under  $D$  is  $\langle c, c' \rangle_D$ .<sup>3</sup> The *average 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  $\gamma$ ,  $\text{SDA}_D(\mathcal{C}, \gamma)$ , is the largest  $d$  such that for all  $\mathcal{C}' \subseteq \mathcal{C}$  with  $|\mathcal{C}'| = d$ ,  $\rho_D(\mathcal{C}') \geq \gamma$ .

*Remark 2.2.5.* For any general and large concept class  $\mathcal{C}$  (such as all one-layer neural nets), we may consider a specific subclass  $\mathcal{C}' \subseteq \mathcal{C}$  and prove lower bounds on learning  $\mathcal{C}'$  in terms of the SDA of  $\mathcal{C}$ . These lower bounds extend to  $\mathcal{C}'$  because if it is hard to learn a subset, then it is hard to learn the whole class.

We will mainly be interested in the statistical dimension in a setting where bounds on pairwise correlations are known. In that case the following lemma holds.

**Lemma 2.2.6** (adapted from [FGR<sup>+</sup>17], Lemma 3.10). *Suppose a concept class  $\mathcal{C}$  has pairwise correlation  $\gamma$ , i.e.  $\langle c, c' \rangle_D \geq \gamma$  for  $c \neq c' \in \mathcal{C}$ , and squared norm at most  $\beta$ , i.e.  $\langle c, c \rangle_D \leq \beta$  for all  $c \in \mathcal{C}$ . Then for any  $\gamma' > 0$ ,  $\text{SDA}_D(\mathcal{C}, \gamma + \gamma') \geq |\mathcal{C}| \frac{\gamma'}{\beta}$ . In particular, if  $\mathcal{C}$  is a class of orthogonal concepts (i.e.  $\gamma = 0$ ) with squared norm bounded by  $\beta$ , then  $\text{SDA}_D(\mathcal{C}, \gamma') \geq |\mathcal{C}| \frac{\gamma'}{\beta}$ .*

*Proof.* Let  $d = |\mathcal{C}| \frac{\gamma'}{\beta}$ , and observe that for any subset  $\mathcal{C}' \subseteq \mathcal{C}$  satisfying  $|\mathcal{C}'| = d$

---

<sup>3</sup>In the  $p$ -concept setting, it is instructive to note that in the notation of [FGR<sup>+</sup>17], this correlation is precisely the distributional correlation  $\chi_{D_0}(D_c, D_{c'})$  of the induced labeled distributions  $D_c$  and  $D_{c'}$  under the reference distribution  $D_0 = D \otimes \text{Unif}^f$ .

$$jCj/d = \frac{\beta \gamma}{\gamma^\theta},$$

$$\begin{aligned} \rho_D(C^\theta) &= \frac{1}{jC^\theta j^2} \sum_{c, c^\theta \in 2C^\theta} jhc, c^\theta i_{Dj} \\ &= \frac{1}{jC^\theta j^2} (jC^\theta j\beta + (jC^\theta j^2 - jC^\theta j)\gamma) \\ &= \gamma + \frac{\beta \gamma}{jC^\theta j} \\ &= \gamma + \gamma^\theta. \end{aligned}$$

□

### 2.3 Orthogonal Family of Neural Networks

We consider neural networks with one hidden layer with activation function  $\phi : \mathbb{R} \rightarrow \mathbb{R}$ , and with one output node that has some activation function  $\psi : \mathbb{R} \rightarrow \mathbb{R}$ . If we take the input dimension to be  $n$  and the number of hidden nodes to be  $m$ , then such a neural network is a function  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  given by

$$f(x) = \psi \left( \sum_{i=1}^m a_i \phi(w_i \cdot x) \right),$$

where  $w_i \in \mathbb{R}^n$  are the weights feeding into the  $i^{\text{th}}$  hidden node, and  $a_i \in \mathbb{R}$  are the weights feeding into the output node. If  $\psi$  takes values in  $[-1, 1]$ , we may also view  $f$  as defining a  $p$ -concept in terms of its conditional mean function.

For our construction, we need our functions to be orthogonal, and we need a lower bound on their norms. For the first property we only need the distribution on the domain to satisfy a relaxed kind of spherical symmetry that we term sign-symmetry, which says that the distribution must look identical on all orthants. To lower bound the norms, we need to assume that the distribution is Gaussian  $N(0, I)$ .

**Assumption 2.3.1** (Sign-symmetry). *For any  $z \in \{-1, 1\}^n$  and  $x \in \mathbb{R}^n$ , let  $x \cdot z$  denote  $(x_1 z_1, \dots, x_n z_n)$ . A distribution  $D$  on  $\mathbb{R}^n$  is sign-symmetric if for any  $z \in \{-1, 1\}^n$  and  $x$  drawn from  $D$ ,  $x$  and  $x \cdot z$  have the same distribution  $D$ .*

**Assumption 2.3.2** (Odd outer activation). *The outer activation  $\psi$  is an odd, increasing function, i.e.  $\psi(-x) = -\psi(x)$ .*

Note that  $\psi$  could be the identity function.

**Assumption 2.3.3** (Inner activation). *The inner activation  $\phi \in L_2(N(0, I))$ .*

The construction of our orthogonal family of neural networks is simple and exploits sign-symmetry.

**Definition 2.3.4** (Family of Orthogonal Neural Networks). Let the domain be  $\mathbb{R}^n$ , let  $\phi : \mathbb{R} \rightarrow \mathbb{R}$  be any well-behaved activation function, and let  $\psi : \mathbb{R} \rightarrow \mathbb{R}$  be any odd function. For an index set  $S \subseteq [n]$ , let  $x_S \in \mathbb{R}^{|S|}$  denote the vector of  $x_i$  for  $i \in S$ . Fix any  $k > 0$ . For any sign-pattern  $z \in \{-1, 1\}^k$ , let  $\chi(z)$  denote the parity  $\prod_i z_i$ . For any index set  $S \subseteq [n]$ , define a one-layer neural network with  $m = 2^k$  hidden nodes,

$$g_S(x) = \sum_{w \in \{-1, 1\}^k} \chi(w) \phi\left(\frac{w \cdot x_S}{k}\right)$$

$$f_S(x) = \psi(g_S(x)).$$

Our orthogonal family is

$$\mathcal{C}_{\text{orth}}(n, k) = \{f_S \mid S \subseteq [n]\}.$$

Notice that the size of this family is  $\binom{n}{k} = n^{\Theta(k)}$  (for appropriate  $k$ ), which is  $n^{\Theta(\log m)}$  in terms of  $m$ . We will take  $k = \Theta(\log n)$ , so that  $m = \text{poly}(n)$  and thus the neural networks are  $\text{poly}(n)$ -sized, and the size of the family is  $n^{\Theta(\log n)}$ , i.e. quasipolynomial in  $n$ .

We now prove that our functions are orthogonal under any sign-symmetric distribution.

**Theorem 2.3.5.** *Let the domain be  $\mathbb{R}^n$ , and let  $D$  be a sign-symmetric distribution on  $\mathbb{R}^n$ . Fix any  $k > 0$ . Then  $\langle f_S, f_T \rangle_D = 0$  for any two distinct  $f_S, f_T \in \mathcal{C}_{\text{orth}}(n, k)$ .*

*Proof.* For the proof, the key property of our construction that we will use is the following: for any sign-pattern  $z \in \{-1, 1\}^n$  and any  $x \in \mathbb{R}^n$ ,

$$f_S(x \cdot z) = \chi_S(z) f_S(x), \quad (2.3.1)$$

where  $\chi_S(z) = \prod_{i \in S} z_i = \chi(z_S)$  is the parity on  $S$  of  $z$ . Indeed, observe first that

$$\begin{aligned} g_S(x \cdot z) &= \sum_{w \in \{-1, 1\}^k} \chi(w) \phi\left(\frac{w \cdot \left(\frac{x \cdot z}{k}\right)_S}{k}\right) \\ &= \sum_{w \in \{-1, 1\}^k} \chi(w) \phi\left(\frac{(w \cdot z_S) \cdot x_S}{k}\right) \\ &= \sum_{w \in \{-1, 1\}^k} \chi(w \cdot z_S) \chi(z_S) \phi\left(\frac{(w \cdot z_S) \cdot x_S}{k}\right) \\ &= \chi(z_S) \sum_{w \in \{-1, 1\}^k} \chi(w) \phi\left(\frac{w \cdot x_S}{k}\right) \quad (\text{replacing } w \cdot z_S \text{ with } w) \\ &= \chi(z_S) g_S(x). \end{aligned}$$

The property then follows since  $\psi$  is odd and  $\psi(av) = a\psi(v)$  for any  $a \in \{-1, 1\}$  and  $v \in \mathbb{R}$ .

Consider  $f_S$  and  $f_T$  for any two distinct  $S, T \subseteq [n]$ . Recall that by the definition of sign-symmetry, for any  $z \in \{-1, 1\}^n$  and  $x$  drawn from  $D$ ,  $x$  and  $x \cdot z$  has the same distribution. Using this and Eq. (2.3.1), we have

$$\begin{aligned} \langle f_S, f_T \rangle_D &= \mathbb{E}_x [f_S(x) f_T(x)] \\ &= \mathbb{E}_z \mathbb{E}_x [f_S(x \cdot z) f_T(x \cdot z)] \quad (\text{sign-symmetry}) \\ &= \mathbb{E}_z \mathbb{E}_x [\chi_S(z) f_S(x) \chi_T(z) f_T(x)] \quad (\text{Eq. (2.3.1)}) \\ &= \mathbb{E}_x \left[ f_S(x) f_T(x) \mathbb{E}_z \chi_S(z) \chi_T(z) \right] \\ &= 0, \end{aligned}$$

since  $\mathbb{E}_z \chi_S(z) \chi_T(z) = 0$  for any two distinct parities  $\chi_S, \chi_T$ .  $\square$

*Remark 2.3.6.* Our proof actually shows that any family of functions satisfying Eq. (2.3.1) is an orthogonal family under any sign-symmetric distribution.

We still need to establish that our functions are nonzero. For this we need to specialize to the Gaussian distribution, as well as consider specific activation functions (a similar analysis can in principle be carried out for other sign-symmetric distributions). For any  $n$  and  $k$ , it follows from Lemma A.1.1 that if the inner activation  $\phi$  has a nonzero Hermite coefficient of degree  $k$  or higher, then the functions in  $\mathcal{C}_{\text{orth}}(n, k)$  are nonzero. The sigmoid, ReLU and sign functions all satisfy this property.

**Corollary 2.3.7.** *Let the domain be  $\mathbb{R}^n$ , and let  $D$  be any sign-symmetric distribution on  $\mathbb{R}^n$ . For any  $\gamma > 0$ ,*

$$\text{SDA}_D(\mathcal{C}_{\text{orth}}(n, k), \gamma) \quad \|\mathcal{C}_{\text{orth}}(n, k)\|_{\gamma} = \binom{n}{k} \gamma.$$

*Here we also assume that all  $c \in \mathcal{C}_{\text{orth}}(n, k)$  are nonzero for our distribution  $D$ .*

*Proof.* Follows from Theorem 2.3.5 and Lemma 2.2.6, using a loose upper bound of 1 on the squared norm.  $\square$

We also need to prove norm lower bounds on our functions for our notions of learning to be meaningful. In Section A.1, we prove the following.

**Theorem 2.3.8.** *Let the inner activation function  $\phi$  be ReLU or sigmoid, and let the outer activation function  $\psi$  be any odd, increasing, continuous function. Let the underlying distribution  $D$  be  $N(0, I_n)$ . Then  $\|f_S\| = \Omega(e^{-\binom{k}{2}})$ , where the hidden constants depend on  $\psi$  and  $\phi$ , for any  $f_S \in \mathcal{C}_{\text{orth}}(n, k)$ .*

With this in hand, we now state our main SQ lower bounds.

**Theorem 2.3.9.** *Let the input dimension be  $n$ , and let the underlying distribution be  $N(0, I_n)$ . Consider  $\mathcal{C}_{\text{orth}}(n, k)$  instantiated with  $\phi = \text{ReLU}$  or sigmoid and  $\psi$  any odd, increasing function (including the identity function), and let  $m = 2^k$  be the hidden*

layer size of each neural net. Let  $A$  be an SQ learner using only inner product queries of tolerance  $\tau$ . For any  $k \geq 2\mathbb{N}$ , there exists  $\tau = 1/n^{\Theta(k)}$  such that  $A$  requires at least  $n^{\Theta(k)}$  queries of tolerance  $\tau$  to learn  $\mathcal{C}_{\text{orth}}(n, k)$  with advantage  $1/\exp(k)$ .

In particular, there exist  $k = \Theta(\log n)$  and  $\tau = 1/n^{\Theta(\log n)}$  such that  $A$  requires at least  $n^{\Theta(\log n)}$  queries of tolerance  $\tau$  to learn  $\mathcal{C}_{\text{orth}}(n, k)$  with advantage  $1/\text{poly}(n)$ . In this case  $m = \text{poly}(n)$ , so that each function in the family has polynomial size. This is our main superpolynomial lower bound.

*Proof.* The proof amounts to careful choices of the parameters  $\epsilon, \gamma$  and  $\tau$  in Corollary 2.3.7 and Corollary 2.4.6. Recall that  $\text{SDA}(\mathcal{C}_{\text{orth}}(n, k), \gamma) \leq n^{\Theta(k)}\gamma$ . We pick  $\gamma = n^{-\Theta(k)}$  appropriately such that  $d = \text{SDA}(\mathcal{C}_{\text{orth}}(n, k), \gamma)$  is still  $n^{\Theta(k)}$ . Theorem 2.3.8 gives us a norm lower bound of  $\exp(-\Theta(k))$ , allowing us to take  $\epsilon = \exp(-\Theta(k))$  and  $\tau = \frac{\epsilon}{\gamma} = n^{-\Theta(k)}$  in Corollary 2.4.6.  $\square$

## 2.4 SQ Lower Bounds

### 2.4.1 SQ Lower Bounds for Real-valued Functions

Prior work [Szö09, Fel12] has already established the following fundamental result, which we phrase in terms of our definition of statistical dimension. For the reader's convenience, we include a proof in Section A.2.

**Theorem 2.4.1.** *Let  $D$  be a distribution on  $X$ , and let  $\mathcal{C}$  be a real-valued concept class over a domain  $X$  such that  $\|c_k c_D\|_D > \epsilon$  for all  $c \in \mathcal{C}$ . Consider any SQ learner that is allowed to make only inner product queries to an SQ oracle for the labeled distribution  $D_c$  for some unknown  $c \in \mathcal{C}$ . Let  $d = \text{SDA}_D(\mathcal{C}, \gamma)$ . Then any such SQ learner needs at least  $\Omega(d)$  queries of tolerance  $\frac{\epsilon}{\gamma}$  to learn  $\mathcal{C}$  up to  $L_2$  error  $\epsilon$ .*

### 2.4.2 SQ Lower Bounds for $p$ -concepts

It turns out to be fruitful to view our learning problem in terms of a decision problem over distributions. We define the problem of distinguishing a valid labeled



distribution from a randomly labeled one, and show a lower bound for this problem. We then show that learning is at least as hard as distinguishing, thereby extending the lower bound to learning as well. Our analysis closely follows that of [FGR<sup>+</sup>17].

**Definition 2.4.2** (Distinguishing between labeled and uniformly random distributions). Let  $\mathcal{C}$  be a class of  $p$ -concepts over a domain  $X$ , and let  $D$  be a distribution on  $X$ . Let  $D_0 = D_{c_0}$  be the randomly labeled distribution  $D \stackrel{\text{Unif}}{\sim} \{g\}$ . Suppose we are given SQ access either to a labeled distribution  $D_c$  for some  $c \in \mathcal{C}$  such that  $c \neq c_0$  or to  $D_0$ . The problem of distinguishing between labeled and uniformly random distributions is to decide which.

*Remark 2.4.3.* Given access to  $D_c$  for some truly boolean concept  $c : X \rightarrow \{0, 1\}$ , it is easy to distinguish any other boolean function  $c^\ell$  from  $c$  since  $\langle c, c^\ell \rangle_D = \frac{1}{2} + \frac{1}{2} \langle c, c^\ell \rangle_D$  (which is information-theoretically optimal as a distinguishing criterion) can be computed using a single inner product query. However, if  $c$  and  $c^\ell$  are  $p$ -concepts,  $\langle c, c \rangle_D$  and  $\langle c^\ell, c^\ell \rangle_D$  are not 1 in general and may be difficult to estimate. It is not obvious how best to distinguish the two, short of directly learning the target.

Considering the distinguishing problem is useful because if we can show that distinguishing itself is hard, then any reasonable notion of learning will be hard as well, including weak learning. We give simple reductions for both our notions of learning.

**Lemma 2.4.4** (Learning is as hard as distinguishing). *Let  $D$  be a distribution over the domain  $X$ , and let  $\mathcal{C}$  be a  $p$ -concept class over  $X$ . Suppose there exists either*

*(a) a weak SQ learner capable of learning  $\mathcal{C}$  up to advantage  $\epsilon$  using  $q$  queries of tolerance  $\tau$ , where  $\tau \leq \epsilon/2$ ; or,*

*(b) an SQ learner capable of learning  $\mathcal{C}$  (assume  $\langle c, c \rangle_D \geq 3\epsilon$  for all  $c \in \mathcal{C}$ ) up to  $L_2$  error  $\epsilon$  using  $q$  queries of tolerance  $\tau$ , where  $\tau \leq \epsilon^2$ . Then there exists a distinguisher that is able to distinguish between an unknown  $D_c$  and  $D_0$  using at most  $q + 1$  queries of tolerance  $\tau$ .*

*Proof.* (a) Run the weak learner to obtain  $\tilde{c}$ . If  $c \neq c_0$ , we know that  $h\tilde{c}, c|_D \geq \epsilon$ , whereas if  $c = c_0$ , then  $h\tilde{c}, c|_D = 0$  no matter what  $\tilde{c}$  is. A single additional query  $(h(x, y) = \tilde{c}(x)y)$  of tolerance  $\epsilon/2$  distinguishes between the two cases.

(b) Run the learner to obtain  $\tilde{c}$ . If  $c \neq c_0$ , i.e.  $kck_D \geq 3\epsilon$ , we know that  $k\tilde{c} - ck_D \geq \epsilon$ , so that by the triangle inequality,  $k\tilde{c}k_D - kck_D \geq k\tilde{c} - ck_D \geq 2\epsilon$ . But if  $c = c_0$ , then  $k\tilde{c}k_D \leq \epsilon$ . An additional query  $(h(x, y) = \tilde{c}(x)^2)$  of tolerance  $\epsilon^2$  suffices to distinguish the two cases.  $\square$

We now prove the main lower bound on distinguishing.

**Theorem 2.4.5.** *Let  $D$  be a distribution over the domain  $X$ , and let  $\mathcal{C}$  be a  $p$ -concept class over  $X$ . Then any SQ algorithm needs at least  $d = \text{SDA}(\mathcal{C}, \gamma)$  queries of tolerance  $\rho_{\gamma}$  to distinguish between  $D_c$  and  $D_0$  for an unknown  $c \in \mathcal{C}$ . (We will consider deterministic SQ algorithms that always succeed, for simplicity.)*

*Proof.* Consider any successful SQ algorithm  $A$ . Consider the adversarial strategy where to every query  $h : X \rightarrow [-1, 1]$  of  $A$  (with tolerance  $\tau = \rho_{\gamma}$ ), we respond with  $\mathbb{E}_{D_0}[h]$ . We can pretend that this is a valid answer with respect to any  $c \in \mathcal{C}$  such that  $|\mathbb{E}_{D_c}[h] - \mathbb{E}_{D_0}[h]| \leq \tau$ . Our argument will be based on showing that each such query rules out fairly few distributions, so that the number of queries required in total is large.

Since we assumed that  $A$  is a deterministic algorithm that always succeeds, it eventually correctly guesses that it is  $D_0$  that it is getting answers from. Say it takes  $q$  queries to do so. For the  $k^{\text{th}}$  query  $h_k$ , let  $S_k$  be the set of concepts in  $\mathcal{C}$  that are ruled out by our response  $\mathbb{E}_{D_0}[h_k]$ :

$$S_k = \{c \in \mathcal{C} \mid |\mathbb{E}_{D_c}[h_k] - \mathbb{E}_{D_0}[h_k]| > \tau\}.$$

We'll show that

(a) on the one hand,  $\bigcup_{k=1}^q S_k = \mathcal{C}$ , so that  $\sum_{k=1}^q |S_k| \geq |\mathcal{C}|$ ,

(b) while on the other,  $|S_k| > |C|/d$  for every  $k$ . Together, this will mean that  $q > d$ .

For the first claim, suppose  $\{S_k\}_{k=1}^q$  were not all of  $\mathcal{C}$ , and indeed say  $c \in \mathcal{C} \setminus \{S_k\}_{k=1}^q$ . This is a distribution that our answers were consistent with throughout, yet one that  $A$ 's solution ( $D_0$ ) is incorrect for. But  $A$  always succeeds, so for it not to have ruled out this  $D_c$  is impossible.

For the second claim, suppose for the sake of contradiction that for some  $k$ ,  $|S_k| > |C|/d$ . By Definition 2.2.4, this means we know that  $\rho_D(S_k) > \gamma$ . One of the key insights in the proof of [Szö09] is that by expressing query expectations entirely in terms of inner products, we gain the ability to apply simple algebraic techniques. To this end, for any query function  $h$ , let  $\hat{h}(x) = (h(x, 1) + h(x, -1))/2$ . Observe that for any  $p$ -concept  $c$ ,

$$\begin{aligned} \langle \hat{h}, c \rangle_{D_c} &= \mathbb{E}_{x \sim D_c} \left[ h(x, 1) \frac{c(x)}{2} \right] - \mathbb{E}_{x \sim D_c} \left[ h(x, -1) \frac{c(x)}{2} \right] \\ &= \mathbb{E}_{x \sim D_c} \left[ h(x, 1) \frac{1 + c(x)}{2} \right] \\ &\quad + \mathbb{E}_{x \sim D_c} \left[ h(x, -1) \frac{1 - c(x)}{2} \right] \\ &= \mathbb{E}_{x \sim D_c} \left[ h(x, 1) \frac{1}{2} \right] - \mathbb{E}_{x \sim D_c} \left[ h(x, -1) \frac{1}{2} \right] \\ &= \mathbb{E}_{D_c} [h] - \mathbb{E}_{D_0} [h], \end{aligned}$$

the difference between the query expectations wrt  $D_c$  and  $D_0$ . Here we have expanded each  $\mathbb{E}_{D_c} [h]$  using the fact that the label for  $x$  is 1 with probability  $(1 + c(x))/2$  and -1 otherwise. Thus  $\langle \hat{h}_k, c \rangle_{D_c}$ , where  $h_k$  is the  $k^{\text{th}}$  query, is greater than  $\tau$  for any  $c \in S_k$ , since  $S_k$  are precisely those concepts ruled out by our response. We will show contradictory upper and lower bounds on the following quantity:

$$\Phi = \left\langle \hat{h}_k, \sum_{c \in S_k} c \cdot \text{sign}(\langle \hat{h}_k, c \rangle_{D_c}) \right\rangle_D.$$

Note that since every query  $h$  satisfies  $k h(\cdot, y) k_D \leq 1$  for all  $y$ , it follows by the triangle inequality that  $k \widehat{h} k_D \leq 1$ . So by Cauchy-Schwarz and our observation that  $\rho_D(S_k) \leq \gamma$ ,

$$\begin{aligned} \Phi^2 &\leq k \widehat{h} k_D^2 \left\| \sum_{c \in \mathcal{C}} c \operatorname{sign}(h \widehat{h}_k, c) \right\|_D^2 \\ &\leq \sum_{c, c' \in \mathcal{C}} j h c, c' i_D j = j S_k j^2 \rho_D(S_k) \leq j S_k j^2 \gamma. \end{aligned}$$

However since  $j h \widehat{h}_k, c i_D j > \tau$ , we also have that  $\Phi = \sum_{c \in \mathcal{C}} j h \widehat{h}_k, c i_D j > j S_k j \tau$ . Since  $\tau = \rho_{\overline{\gamma}}$ , this contradicts our upper bound and in turn completes the proof of our second claim. And as noted earlier, the two claims together imply that  $q \leq d$ .  $\square$

The final lower bounds on learning thus obtained are stated as a corollary for convenience. The proof follows directly from Lemma 2.4.4 and Theorem 2.4.5.

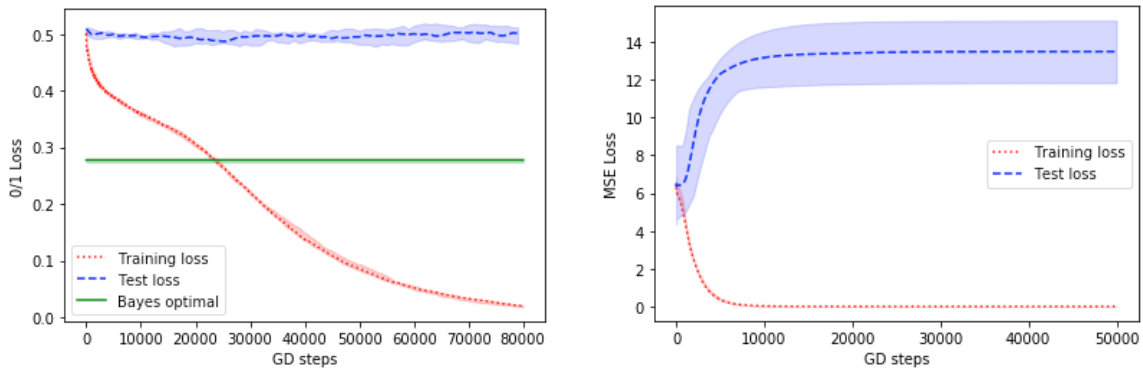
**Corollary 2.4.6.** *Let  $D$  be a distribution over the domain  $X$ , and let  $\mathcal{C}$  be a  $p$ -concept class over  $X$ . Let  $\gamma, \tau$  be such that  $\rho_{\overline{\gamma}} \leq \tau$ . Let  $d = \text{SDA}(\mathcal{C}, \gamma)$ .*

(a) *Let  $\epsilon$  be such that  $\tau \leq \epsilon^2$ , and assume  $k c k_D \leq 3\epsilon$  for all  $c \in \mathcal{C}$ . Then any SQ learner learning  $\mathcal{C}$  up to  $L_2$  error  $\epsilon$  requires at least  $d - 1$  queries of tolerance  $\tau$ .*

(b) *Let  $\epsilon$  be such that  $\tau \leq \epsilon/2$ . Then any weak SQ learner learning  $\mathcal{C}$  up to advantage  $\epsilon$  requires at least  $d - 1$  queries of tolerance  $\tau$ .*

## 2.5 Experiments

We include experiments for both regression and classification. We train an overparameterized neural network on data from our function class, using gradient descent. We find that we are able to achieve close to zero training error, while test error remains high. This is consistent with our lower bound for these classes of functions.



(a) Learning a softmax of a one-layer tanh network      (b) Learning a linear combination of tanhs network

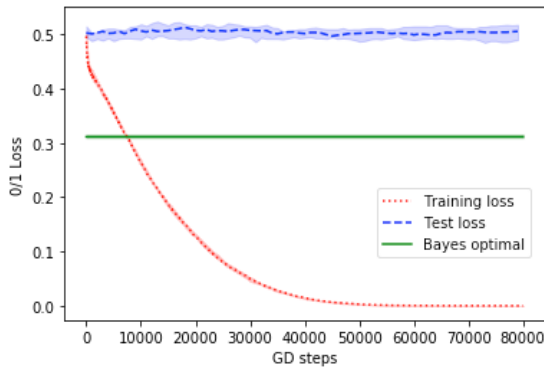
Figure 2.1: In (a) the target function is a softmax ( 1 labels) of a sum of  $2^7$  tanh activations with  $n = 14$ ; in (b) the labels are obtained similarly but without the softmax. In both cases, we train a 1-layer neural network with  $5 \cdot 2^7 = 640$  tanh units (hence 10241 parameters) using a training set of size 6000 and a test set of size 1000, with the learning rate set to 0.01. For (a) we take the sign of this trained network and measure its training and testing 0/1 loss; for (b) we measure the train and test square-loss of the learned network directly. In (a) we also plot the test error of the bayes optimal network (sign of the target function).

For classification, we use a training set of size  $T$  of data corresponding to  $f \in \mathcal{C}_{\text{orth}}(n, k)$  instantiated with  $\phi = \tanh$  and  $\psi = \tanh$ . We draw  $x \sim \mathcal{N}(0, I_n)$ . For each  $x$ ,  $y$  is picked randomly from  $f^{-1}g$  in such a way that  $\mathbb{E}[y/x] = f(x)$ . Since the outer activation  $\psi$  is  $\tanh$ , this can be thought of as applying a softmax to the network’s output, or as the Boolean label corresponding to a logit output. We train a sum of tanh network (i.e. a network in which the inner activation is  $\tanh$  and no outer activation is applied) on this data using gradient descent on squared loss, threshold the output, and plot the resulting 0/1 loss. See Fig. 2.1(a). This setup models a common way in which neural networks are trained for classification problems in practice.

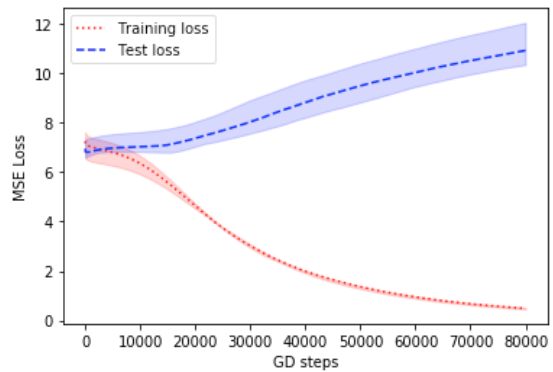
For regression, we use a training set of size  $T$  of data corresponding to  $f \in \mathcal{C}_{\text{orth}}(n, k)$  instantiated with  $\phi = \tanh$  and  $\psi$  being the identity. We draw  $x \sim \mathcal{N}(0, I_n)$ , and  $y = f(x)$ . We train a sum of tanh network on this data using gradient descent on squared loss, which we plot in Fig. 2.1(b). This setup models the natural way of using neural networks for regression problems.

In both cases, we train neural networks whose number of parameters considerably exceeds the amount of training data. In all our experiments, we plot the median over 10 trials and shade the inter-quartile range of the data.

Similar results hold with the inner activation  $\phi$  being ReLU instead of  $\tanh$ , and are shown in Fig. 2.2.



(a) Learning a softmax of a one-layer ReLU network



(b) Learning a linear combination of ReLUs

Figure 2.2: In (a) the target function is a softmax ( 1 labels) of a sum of  $2^8$  ReLU activations with  $n = 14$ ; in (b) the labels are obtained similarly but without the softmax. In both cases, we train a 1-layer neural network with  $5 \cdot 2^8 = 1280$  ReLU units (hence 20481 parameters) using a training set of size 6000 and a test set of size 1000, with the learning rate set to 0.005 for classification and 0.002 for regression. For (a) we take the sign of this trained network and measure its training and testing 0/1 loss; for (b) we measure the train and test square loss of the learned network directly. In (a) we also plot the test error of the bayes optimal network (sign of the target function).

# Chapter 3: Hardness of Noise-Free Learning for Two-Hidden-Layer Neural Networks

## 3.1 Introduction

In this chapter we extend a central line of research proving representation-independent hardness results for learning classes of neural networks. We will consider arguably the simplest possible setting: given samples  $(x_1, y_1), \dots, (x_n, y_n)$  where for every  $i \in [n]$ ,  $x_i$  is sampled independently from some distribution  $D$  over  $\mathbb{R}^d$  and  $y_i = f(x_i)$  for an unknown neural network  $f : \mathbb{R}^d \rightarrow \mathbb{R}$ , the goal is to output any function  $\hat{f}$  for which  $\mathbb{E}_{x \sim D}[(f(x) - \hat{f}(x))^2]$  is small. This model is often referred to as the *realizable* or *noise-free* setting.

This problem has long been known to be computationally hard for discrete input distributions. For example, if  $D$  is supported over a discrete domain like the Boolean hypercube, then we have a variety of hardness results based on cryptographic/average-case assumptions [KS09, DLSS14, Dan16a, DV20b, DV21].

Over the last few years there has been a very active line of research on the complexity of learning with respect to continuous distributions, the most widely studied case being the assumption that  $D$  is a standard Gaussian in  $d$  dimensions. A rich algorithmic toolbox has been developed for the Gaussian setting [JSA15, ZSJ+17b, BG17b, LY17, Tia17, GKM18, GLM18b, BJW19, ZYWG19b, DGK+20, LMZ20, DK20, ATV21, CKM20, SZB21, VSS+22], but all known efficient algorithms can only handle networks with a single hidden layer, that is, functions of the form  $f(x) = \sum_{i=1}^k \lambda_i \sigma(\langle w_i, x \rangle)$ . This motivates the following well-studied question:

*Are there fundamental barriers to learning neural networks with two hidden layers?*

(3.1.1)

Two distinct lines of research, one using cryptography and one using the statistical query (SQ) model, have made progress towards solving this question.



In the cryptographic setting, [DV21] showed that the existence of a certain class of pseudorandom generators, specifically local pseudorandom generators with polynomial stretch, implies superpolynomial lower bounds for learning ReLU networks with *three* hidden layers.

For SQ learning, work of [GGJ<sup>+</sup>20] and [DKKZ20] gave the first superpolynomial *correlational* SQ (CSQ) lower bounds for learning even one-hidden-layer neural networks. Notably, however, there are strong separations between SQ and CSQ [APVZ14, ADHV19, CKM20], and the question of whether a general SQ algorithm exists remained an interesting open problem. In fact, Vempala and Wilmes [VW19a] showed that general SQ lower bounds might be impossible to achieve for learning real-valued neural networks. For any family of networks satisfying a simple non-degeneracy condition (see Section 3.1.1), they gave an algorithm that succeeded using only polynomially many statistical queries. As such, the prevailing conventional wisdom was that noise was required in the model to obtain full SQ lower bounds.

The main contribution of this chapter is to answer Question 3.1.1 by giving both general SQ lower bounds and cryptographic hardness results (based on the Learning with Rounding or LWR assumption) for learning ReLU networks with two hidden layers and polynomially bounded weights.<sup>1</sup> We note that our SQ lower bound is the first of its kind for learning ReLU networks of *any* depth. We also show how to extend our results to the setting where the learner has label query access to the unknown network.

**SQ Lower Bound** We state an informal version of our main SQ lower bound:

**Theorem 3.1.1** (Full SQ lower bound for two hidden layers (informal), see Theorem 3.4.1). *Any SQ algorithm for learning  $\text{poly}(d)$ -sized two-hidden-layer ReLU net-*

---

<sup>1</sup>Note that if the weights were allowed to be arbitrarily large, it is well-known to be trivial to obtain hardness over Gaussian inputs from hardness over Boolean inputs: simply approximate the sign function arbitrarily well and convert all but an arbitrarily small fraction of Gaussian inputs to bitstrings.

| Reference                     | Num. hidden layers | Model of hardness                                   |
|-------------------------------|--------------------|---|
| [DKKZ20, GGJ <sup>+</sup> 20] | 1                  | Correlational SQ                                    |
| [DV21]                        | 3                  | Cryptographic<br>(assuming existence of local PRGs) |
| <i>This work</i>              | 2                  | Full SQ   |
| <i>This work</i>              | 2                  | Cryptographic<br>(assuming hardness of LWR)         |

Table 3.1: Summary of known and new superpolynomial lower bounds for learning noise-free shallow ReLU networks over Gaussian inputs up to sufficiently small (but non-negligible) error. (Definitions and terminology may be found in Section 3.2.)

*works over  $N(0, \text{Id}_d)$  up to squared loss  $1/\text{poly}(d)$  must use at least  $d^{\omega(1)}$  queries, or have query tolerance that is negligible in  $d$ .*

We stress that this bound holds unconditionally, independent of any cryptographic assumptions. This simultaneously closes the gap between the hardness result of [DV21] and the positive results on one-hidden-layer networks [JSA15, ZSJ<sup>+</sup>17b, GLM18b, ATV21, DK20] and goes against the conventional wisdom that one cannot hope to prove full SQ lower bounds for learning real-valued functions in the realizable setting.

We also note that unlike previous CSQ lower bounds which are based on orthogonal function families and crucially exploit cancellations specific to the Gaussian distribution, our Theorem 3.1.1 and other hardness results in this chapter easily extend to any reasonably anticoncentrated and symmetric product distribution over  $\mathbb{R}^d$ ; see Remark 3.3.11.

**Cryptographic Lower Bound** While Theorem 3.1.1 rules out almost all known approaches for provably learning neural networks (e.g. method of moments/tensor decomposition [JSA15, ZSJ<sup>+</sup>17b, GLM18b, BJW19, DGK<sup>+</sup>20, DK20, ATV21], noisy gradient descent [BG17b, LY17, Tia17, GKM18, ZYWG19b, LMZ20], and filtered

PCA [CKM20]), it does not preclude the existence of a non-SQ algorithm for doing so. Indeed, a number of recent works [BRST21, SZB21, ZSWB22, DK21] have ported algorithmic techniques like lattice basis reduction [LLL82], traditionally studied in the context discrete settings like cryptanalysis, to learning problems over continuous domains for which there is no corresponding SQ algorithm.

Our next result shows however that under a certain cryptographic assumption, namely hardness of *Learning with Rounding (LWR) with polynomial modulus* [BPR12, AKPW13, BGM<sup>+</sup>16] (see Section 3.2), no polynomial-time algorithm can learn two-hidden-layer neural networks from Gaussian examples.

**Theorem 3.1.2** (Cryptographic hardness result (informal), see Theorem 3.5.1). *Suppose there exists a  $\text{poly}(d)$ -time algorithm for learning  $\text{poly}(d)$ -sized two-hidden-layer ReLU networks over  $N(0, \text{Id}_d)$  up to squared loss  $1/\text{poly}(d)$ . Then there exists a quasipolynomial-time algorithm for LWR with polynomial modulus.*

Note that here we may actually improve the LWR hardness assumption required from quasipolynomial to any mildly superpolynomial function of the security parameter (see Remark 3.5.2).

Under LWR with polynomial modulus, we also show the first hardness result for learning *one hidden layer* ReLU networks over the uniform distribution on  $f_0, 1g^d$  (see Theorem 3.5.3).

In Section 3.2, we discuss existing hardness evidence for LWR as well as its relation to more standard assumptions like Learning with Errors. From a negative perspective, Theorem 3.1.2 suggests that the aforementioned lattice-based algorithms for continuous domains are unlikely to yield new learning algorithms for two-hidden-layer networks, because even their more widely studied *discrete* counterparts have yet to break LWR. From a positive perspective, in light of the prominent role LWR and its variants have played in a number of practical proposals for post-quantum cryptography [CKLS18, BGML<sup>+</sup>18, JZ16, DKRV18], Theorem 3.1.2 offers a new avenue for stress-testing these schemes.

**Query Learning Lower Bound** One additional benefit of our techniques is that they are flexible enough to accommodate other learning models beyond traditional PAC learning. To illustrate this, for our final result we show hardness of learning neural networks from *label queries*. In this setting, the learner is much more powerful: rather than sample or SQ access, they are given the ability to query the value  $f(x)$  of the unknown function  $f$  at any desired point  $x$  in  $\mathbb{R}^d$ , and the goal is still to output a function  $\hat{f}$  for which  $\mathbb{E}[(f(x) - \hat{f}(x))^2]$  is small. The expectation here is with respect to some specified distribution, which we will take to be  $\mathcal{N}(0, \text{Id}_d)$ , though as before, our techniques will apply to any reasonably anticoncentrated, symmetric product distribution over  $\mathbb{R}^d$ .

In recent years, this question has received renewed interest from the security and privacy communities in light of *model extraction attacks*, which attempt to reverse-engineer neural networks found in publicly deployed systems [TJ<sup>+</sup>16, MSDH19, PMG<sup>+</sup>17, JCB<sup>+</sup>20, RK20, JWZ20, DG21]. Recent work [CKM21] has shown that in this model, there is an efficient algorithm for learning arbitrary one-hidden-layer ReLU networks that is truly polynomial in all relevant parameters. We show that under plausible cryptographic assumptions about the existence of simple pseudorandom function (PRF) families (see Section 3.6) which may themselves be based on standard number theoretic or lattice-based cryptographic assumptions, such a guarantee is impossible for general *constant-depth* ReLU networks.

**Theorem 3.1.3** (Label query hardness (informal), see Theorem 3.6.1). *If either the decisional Diffie-Hellman or the Learning with Errors assumption holds, then the class of  $\text{poly}(d)$ -sized constant-depth ReLU networks from  $\mathbb{R}^d$  to  $\mathbb{R}$  is not learnable up to small constant squared loss  $\epsilon$  over  $\mathcal{N}(0, \text{Id}_d)$  even using label queries over all of  $\mathbb{R}^d$ .*

Note that the connection between PRFs and hardness of learning from label queries over *discrete domains* is a well-known connection dating back to Valiant [Val84].

To our knowledge, however, Theorem 3.1.3 is the first hardness result for query learning over continuous domains.

### 3.1.1 Discussion and Related Work

**Hardness for learning neural networks.** There are a number of works [BR89, Vu06, KS09, LSSS14, GKKT17a, DV20b] showing hardness for *distribution-free* learning of various classes of neural networks.

As for hardness of distribution-specific learning, several works have established lower bounds with respect to the Gaussian distribution. Apart from the works [GGJ<sup>+</sup>20, DKKZ20, DV21] from the introduction which are most closely related to the present work, we also mention the works of [KK14a, GKK19, GGK20, DKZ20a] which showed hardness for *agnostically* learning halfspaces and ReLUs, [Sha18b] which showed hardness for learning periodic activations with gradient-based methods, [SVWX17] which showed lower bounds against SQ algorithms for learning one-hidden-layer networks using *Lipschitz* statistical queries and large tolerance, and [SZB21] which showed lattice-based hardness of learning one-hidden-layer networks when the labels  $y_i$  have been perturbed by bounded *adversarially chosen* noise. Our approach has similarities to the “Gaussian lift” as studied by Klivans and Kothari [KK14a]. Their approach, however, required noise in the labels, whereas we are interested in hardness in the strictly *realizable setting*. We also remark that [DGKP20, AAK21] showed *correlational* SQ lower bounds for learning random depth- $\omega(\log n)$  neural networks over *Boolean inputs* which are uniform over a halfspace.

There have also been works on hardness of learning from label queries over *discrete domains* and for more “classical” concept classes like Boolean circuits [Fel09, CGV15, Val84, Kha95, AK95].

Lastly, we remark on how our results relate to [CKM20], which gives the only known upper bound for learning neural networks over Gaussian inputs beyond one hidden layer. They showed that learning ReLU networks of arbitrary depth is “fixed-

parameter tractable” in the sense that there is a fixed function  $g(k, \epsilon)$  in the size  $k$  of the network and target error  $\epsilon$  for which the time complexity is at most  $g(k, \epsilon) \text{ poly}(d)$ , and their algorithm can be implemented in SQ. That said, this does not contradict our lower bounds for two reasons: 1) their algorithm only applies to networks without biases, 2) in our lower bound constructions,  $k$  scales polynomially in  $d$ .

**SQ lower bounds for real-valued functions.** A recurring conundrum in the literature on SQ lower bounds for supervised learning has been whether one can show SQ hardness for learning *real-valued* functions. SQ lower bounds for Boolean functions are typically shown by lower bounding the *statistical dimension* of the function class, which essentially corresponds to the largest possible set of functions in the class which are all approximately pairwise orthogonal. Indeed, the content of the hardness results of [GGJ<sup>+</sup>20, DKKZ20] was to prove lower bounds on the statistical dimension of one-hidden-layer networks. Unfortunately, for real-valued functions, statistical dimension lower bounds only imply CSQ lower bounds. As discussed in [GGJ<sup>+</sup>20], the class of  $d$ -variate Hermite polynomials of degree- $\ell$  is pairwise orthogonal and of size  $d^{O(\ell)}$ , which translates to a CSQ lower bound of  $d^{(\ell)}$ . Yet there exist SQ algorithms for learning Hermite polynomials in far fewer queries [APVZ14, ADHV19].

Further justification for the difficulty of proving SQ lower bounds for real-valued functions came from [VW19a], which observed that for any real-valued learning problem satisfying a seemingly innocuous non-degeneracy assumption—namely that for any pair of functions  $f, g$  in the class, the probability under the input distribution  $D$  that  $f(x) = g(x)$  is zero—there is an efficient “cheating” SQ algorithm (see Proposition 4.1 therein). The SQ lower bound shown in the present work circumvents this proof barrier by exhibiting a family of neural networks for which any pair of networks agrees on a set of inputs with Gaussian measure *bounded away from zero*.

**Open questions** While our results settle Question 3.1.1, a number of intriguing gaps between our lower bounds and existing upper bounds remain open:

- **General one-hidden-layer networks.** Despite the considerable amount of work on learning one-hidden-layer networks over Gaussian inputs, all known positive results that run in polynomial time in all parameters (input dimension  $d$ , network size  $k$ , inverse error  $1/\epsilon$ ) still need to make various assumptions on the structure of the network. Remarkably, it is even open whether one-hidden-layer ReLU networks with *positive output layer weights* (i.e. “sums of ReLUs”) can be learned in polynomial time, the best known guarantee being the  $(k/\epsilon)^{\log^2 k}$   $\text{poly}(d/\epsilon)$ -time algorithm of [DK20]. As for general one-hidden-layer ReLU networks, it is still open whether they can even be learned in time  $d^{O(k)}$   $\text{poly}(1/\epsilon)$ , the best known guarantee being the  $k^{\text{poly}(k/\epsilon)}$   $\text{poly}(d)$ -time algorithm of [CKM20].
- **Query learning shallow networks.** While Theorem 3.1.3 establishes that above a certain constant depth, ReLU networks cannot be learned even from label queries over the Gaussian distribution. It would be interesting to close the gap between this and the positive result of [CKM21] which only applies to one-hidden-layer networks, although fully settling this seems closely related to the question of what are the shallowest possible Boolean circuits needed to implement pseudorandom functions, a longstanding open question in circuit complexity.

### 3.1.2 Technical Overview

Our work will build on a recent approach of Daniely and Vardi [DV21], who developed a simple and clever technique for lifting discrete functions to the Gaussian domain entirely in the realizable setting. Our main contributions are to (1) make their lifting procedure more efficient so that two hidden layers suffice and (2) show how to apply the lift in a variety of models beyond PAC. For the purposes of this overview we will take the domain of our discrete functions to be  $\{0, 1\}^d$ , but our techniques extend to  $\mathbb{Z}_q^d$  with  $q = \text{poly}(d)$ .

**Daniely–Vardi (DV) lift.** At a high level, the DV lift is a transformation mapping a Boolean example  $(x, y)$  labeled by a hard-to-learn Boolean function  $f$  to a Gaussian example  $(z, \tilde{y})$  labeled by a (real-valued) ReLU network  $f^{\text{DV}}$  that behaves similarly to  $f$  in that  $f^{\text{DV}}(z)$  approximates  $f(\text{sign}(z))$ , where for us  $\text{sign}(t)$  denotes  $\mathbb{1}[t > 0]$  and is applied elementwise. The key idea is to use a continuous approximation  $\widetilde{\text{sign}}$  of the sign function, and to pair it with a “soft indicator” function  $\text{bad} : \mathbb{R}^d \rightarrow \mathbb{R}_+$  that is large whenever  $\text{sign}(z) \neq \widetilde{\text{sign}}(z)$ , and that can be implemented as a one-hidden-layer network independent of the target function. One can show that whenever  $f$  is realizable as an  $L$ -hidden-layer network over  $\{0, 1\}^d$ , the function  $f^{\text{DV}}(z) = \text{ReLU}(f(\widetilde{\text{sign}}(z)) - \text{bad}(z))$  can be implemented as an  $(L + 2)$ -hidden-layer network satisfying

$$f^{\text{DV}}(z) = \text{ReLU}(f(\text{sign}(z)) - \text{bad}(z)). \quad (3.1.2)$$

This property allows us to generate synthetic Gaussian labeled examples  $(z, f^{\text{DV}}(z))$  from Boolean labeled examples  $(x, f(x))$ , and thereby reduce the problem of learning  $f$  to that of learning  $f^{\text{DV}}$ . For a fuller overview, see Section 3.3.1.

**Improving the DV lift.** Our first technical contribution is to introduce a more efficient lift which only requires one extra hidden layer. Our starting point is to observe that a variety of hard-to-learn Boolean functions  $f$  like parity and LWR take the form  $f(x) = \sigma(h(x))$  for some ReLU network  $h$  whose range  $T$  over Boolean inputs is a discrete subset of  $[0, \text{poly}(d)]$  of *polynomially bounded* size, and for some function  $\sigma : T \rightarrow [0, 1]$ . For such *compressible* functions (see Definition 3.3.1), one can write  $f(x) = \sigma(h(x)) = \sum_{t \in T} \sigma(t) \mathbb{1}[h(x) = t]$ . Again, we would like to implement lifted function  $f^{\text{M}} : \mathbb{R}^d \rightarrow \mathbb{R}$  using  $\widetilde{\text{sign}}$  and  $\text{bad}$  so that it approximates  $f(\text{sign}(z))$  except when  $\text{bad}$  indicates that  $\widetilde{\text{sign}} \neq \text{sign}$ . To this end, we might hope to implement, say,

$$f^{\text{M}}(z) = \sum_{t \in T} \sigma(t) \mathbb{1}[h(\widetilde{\text{sign}}(z)) = t] \mathbb{1}[\text{bad}(z_j) = 1].$$

Here we now view  $\text{bad}$  as a univariate function, and whenever it is small, we can be sure  $\widetilde{\text{sign}} = \text{sign}$ . Suppose that we could build a one-hidden-layer network  $N(s_1, \dots, s_d; t)$



that behaves like  $\mathbb{1}[t = 0]\mathbb{1}[\mathcal{G}_j : s_j = 1]$ . Then we could realize  $f^M$  as an  $(L + 1)$ -hidden-layer network:

$$f^M(z) = \sum_{t \in \mathcal{Z}^T} \sigma(t) N(\text{bad}(z_1), \dots, \text{bad}(z_d); h(\widetilde{\text{sign}}(z)) \cdot t). \quad (3.1.3)$$

While many natural attempts to build such an  $N$  run into difficulties, we construct a suitably relaxed version of  $N$  that turns out to suffice for the reduction. To gain some intuition for our construction, the starting observation is that the following inclusion-exclusion type formula vanishes identically whenever any of the  $s_j$  is 1:

$$\psi(s_1, s_2, s_3) - \psi(1, s_2, s_3) - \psi(s_1, 1, s_3) - \psi(s_1, s_2, 1) \quad (3.1.4)$$

$$+ \psi(s_1, 1, 1) + \psi(1, s_2, 1) + \psi(s_1, 1, 1) - \psi(1, 1, 1). \quad (3.1.5)$$

For a suitable choice of  $\psi$ , one might hope to build  $N$  out of such a formula by taking  $s_j = \text{bad}(z_j)$  for every  $j$ . But the natural generalization of this expression to  $d$  inputs would have size  $2^d$ , which runs the risk of rendering the resulting SQ lower bounds vacuous. Our final construction (Lemma 3.3.10) instead resembles a truncated inclusion-exclusion type formula of only quasipolynomial size, which may be of independent interest. Since the SQ lower bounds for Boolean functions that we build on are exponential, by a simple padding argument we still obtain a superpolynomial SQ lower bound for our lifted functions.

**Hard one-hidden-layer Boolean functions and LWR.** To use this lift for Theorems 3.1.1 and 3.1.2, we need one-hidden-layer networks that are *compressible* and hard to learn over uniform Boolean inputs. For SQ lower bounds, we can simply start from parities, for which there are exponential SQ lower bounds, and which turn out to be easily implementable by compressible one-hidden-layer networks. For cryptographic hardness, Daniely and Vardi [DV21] used certain one-hidden-layer Boolean networks that arise from the cryptographic assumption that local PRGs exist (see Section A.4.1 therein). Unfortunately, these functions are not compressible. For

this reason, we work instead with LWR: it turns out that the LWR functions are compressible and, conveniently, the hardness assumption directly involves uniform discrete inputs.

**Hardness beyond PAC.** While the DV lift is *a priori* only for showing hardness of example-based PAC learning, we can extend it to the SQ and label query models by simple simulation arguments.

## 3.2 Preliminaries

### 3.2.1 Notation

We use  $\text{Unif}(S)$  to denote the uniform distribution over a set  $S$ . We use  $U_d$  as shorthand for  $\text{Unif}(\mathbb{F}_0, 1\mathbb{G}^d)$ . We use  $\mathcal{N}(0, \text{Id}_d)$  (or sometimes  $\mathcal{N}_d$  for short) to denote the standard Gaussian, and  $j\mathcal{N}(0, \text{Id}_d)j$  (or  $j\mathcal{N}_d j$  for short) to denote the positive standard half-Gaussian (i.e.,  $g = j\mathcal{N}(0, \text{Id}_d)j$  if  $g = |z|$  for  $z \sim \mathcal{N}(0, \text{Id}_d)$ ). We use  $[n]$  to denote  $\mathbb{F}_1, \dots, n\mathbb{G}$ .

For  $q > 0$ ,  $Z_q$  will denote the integers modulo  $q$ , which we will identify with  $\mathbb{F}_0, \dots, (q-1)\mathbb{G}$ . We use  $Z_q/q$  to denote  $\mathbb{F}_0, 1/q, \dots, (q-1)/q\mathbb{G}$ . Our discrete functions will in general have domain  $Z_q^d$  for some  $q$ . The  $q = 2$  case, namely Boolean functions, have domain  $\mathbb{F}_0, 1\mathbb{G}^d$ . For the purposes of this chapter,  $\text{sign} : \mathbb{R} \rightarrow \mathbb{F}_0, 1\mathbb{G}$  is defined as  $\text{sign}(t) = \mathbb{1}[t > 0]$ . We will extend this to  $Z_q$  by defining  $\text{thres}_q : \mathbb{R} \rightarrow Z_q$  in terms of a certain partition of  $\mathbb{R}$  into  $q$  intervals  $I_0, \dots, I_{q-1}$  (formally defined later) as the piecewise constant function that takes on value  $k$  on  $I_k$  for each  $k \in Z_q$ . Scalar functions and scalar arithmetic applied to vectors act elementwise. We say a quantity is *negligible* in a parameter  $n$ , denoted  $\text{negl}(n)$ , if it decays as  $1/n^{\omega(1)}$ .

A *one-hidden-layer* ReLU network mapping  $\mathbb{R}^d$  to  $\mathbb{R}$  is a linear combination of ReLUs, that is, a function of the form

$$F(x) = W_1 \text{ReLU}(W_0 x + b_0) + b_1, \tag{3.2.1}$$

where  $W_0 \in \mathbb{R}^{k \times d}$ ,  $W_1 \in \mathbb{R}^{1 \times k}$ ,  $b_0 \in \mathbb{R}^k$ , and  $b_1 \in \mathbb{R}$ . A *two-hidden-layer* ReLU network mapping  $\mathbb{R}^d$  to  $\mathbb{R}$  is a linear combination of ReLUs of one-hidden-layer networks, that is, a function of the form

$$F(x) = W_2 \text{ReLU} (W_1 \text{ReLU} (W_0 x + b_0) + b_1) + b_2, \quad (3.2.2)$$

where  $W_0 \in \mathbb{R}^{k_0 \times d}$ ,  $W_1 \in \mathbb{R}^{k_1 \times k_0}$ ,  $W_2 \in \mathbb{R}^{1 \times k_1}$ ,  $b_0 \in \mathbb{R}^{k_0}$ ,  $b_1 \in \mathbb{R}^{k_1}$ , and  $b_2 \in \mathbb{R}$ . Our usage of the term *hidden layer* thus corresponds to a *nonlinear* layer.

### 3.2.2 Learning models

Let  $\mathcal{C}$  be a function class mapping  $\mathbb{R}^d$  to  $\mathbb{R}$ , and let  $D$  be a distribution on  $\mathbb{R}^d$ . We consider various learning models where the learner is given access in different ways to labeled data  $(x, f(x))$  for an unknown  $f \in \mathcal{C}$  and must output a (possibly randomized) predictor that achieves (say) squared loss  $\epsilon$  for any desired  $\epsilon > 0$ . In the traditional PAC model, access to the data is in the form of iid labeled examples  $(x, f(x))$  where  $x \sim D$ , and the learner is considered efficient if it succeeds using  $\text{poly}(d, 1/\epsilon)$  time and sample complexity. In the Statistical Query (SQ) model [Kea98, Rey20], access to the data is through an SQ oracle. Given a bounded query  $\phi : \mathbb{R}^d \rightarrow \mathbb{R} \in [-1, 1]$  and a tolerance  $\tau > 0$ , the oracle may respond with any value  $v$  such that  $|v - \mathbb{E}_{x \sim D}[\phi(x, f(x))]| \leq \tau$ . A correlational query is one that is linear in  $y$ , i.e. of the form  $\phi(x, y) = \tilde{\phi}(x)y$  for some  $\tilde{\phi}$ , and a correlational SQ (CSQ) learner is one that is only allowed to make CSQs. An SQ learner is considered efficient if it succeeds using  $\text{poly}(d, 1/\epsilon)$  queries and tolerance  $\tau = 1/\text{poly}(d, 1/\epsilon)$ . Finally, in the label query model, the learner is allowed to request the value of  $f(x)$  for any desired  $x$ , and is considered efficient if it succeeds using  $\text{poly}(d, 1/\epsilon)$  time and queries.

### 3.2.3 Learning with Rounding

The Learning with Rounding (LWR) problem [BPR12] is a close cousin of the well-known Learning with Errors (LWE) problem [Reg09], except with deterministic rounding in place of random additive errors.

**Definition 3.2.1.** Fix moduli  $p, q \geq 2 \in \mathbb{N}$ , where  $p < q$ , and let  $n$  be the security parameter. For any  $w \in \mathbb{Z}_q^n$ , define  $f_w : \mathbb{Z}_q^n \rightarrow \mathbb{Z}_p/p$  by

$$f_w(x) = \frac{1}{p}bw \quad x e_p = \frac{1}{p}b \frac{p}{q}(w \cdot x \bmod q)e,$$

where  $\lfloor t \rfloor$  is the closest integer to  $t$ . In the  $\text{LWR}_{n,p,q}$  problem, the secret  $w$  is drawn randomly from  $\mathbb{Z}_q^n$ , and we must distinguish between labeled examples  $(x, y)$  where  $x \in \mathbb{Z}_q^n$  and either  $y = f_w(x)$  or  $y$  is drawn independently from  $\text{Unif}(\mathbb{Z}_p/p)$ . The  $\text{LWE}_{n,q,B}$  problem is similar, except that  $y \in \mathbb{Z}_q/q$  is either  $\frac{1}{q}((w \cdot x + e) \bmod q)$  for some  $e \in \mathbb{Z}_q$  sampled from a carefully chosen distribution, e.g. discrete Gaussian, such that  $\Pr[e \in B] = \epsilon$  except with  $\text{negl}(n)$  probability, or is drawn from  $\text{Unif}(\mathbb{Z}_q/q)$ .

*Remark 3.2.2.* Traditionally the LWR problem is stated with labels lying in  $\mathbb{Z}_p$  instead of  $\mathbb{Z}_p/p$ , although both are equivalent since the moduli  $p, q$  may be assumed to be known to the learner. The choice of  $\mathbb{Z}_p/p$  is simply a convenient way to normalize labels to lie in  $[0, 1]$ . For consistency, we similarly normalize LWE labels to lie in  $\mathbb{Z}_q/q$ .

It is known that  $\text{LWE}_{n,q,B}$  is as hard as worst-case lattice problems when  $q = \text{poly}(n)$  and  $B = q/\text{poly}(n)$  (see e.g. [Reg10, Pei16] for surveys). Yet this is not known to directly imply the hardness of  $\text{LWR}_{n,p,q}$  in the regime in which  $p, q$  are both  $\text{poly}(n)$ , which is the one we will be interested in as  $p, q$  will dictate the size of the hard networks that we construct in the proof of our cryptographic lower bound.

Unfortunately, in this polynomial modulus regime, it is only known how to reduce from LWE to LWR *when the number of samples is bounded relative to the modulus* [AKPW13, BGM<sup>+</sup>16]. For instance, the best known reduction in this regime obtains the following hardness guarantee:

**Theorem 3.2.3** ([BGM<sup>+</sup>16]). *Let  $n$  be the security parameter, let  $p, q \geq 1$  be moduli, and let  $m, B \geq 0$ . Assuming  $q \geq \Omega(mBp)$ , any distinguisher capable of solving  $\text{LWR}_{n,p,q}$  using  $m$  samples implies an efficient algorithm for  $\text{LWE}_{n,q,B}$ .*

For our purposes, Theorem 3.2.3 is not enough to let us base our Theorem 3.1.2 off of LWE, as we are interested in the regime where the learner has an *arbitrary* polynomial number of samples.

LWR with polynomial modulus and arbitrary polynomial samples is nevertheless conjectured to be as hard as worst-case lattice problems [BPR12] and has already formed the basis for a number of post-quantum cryptographic proposals [DKRV18, CKLS18, BGML<sup>+</sup>18, JZ16]. We remark that one piece of evidence in favor of this conjecture is a reduction from a less standard variant of LWE in which the usual discrete Gaussian errors are replaced by errors uniformly sampled from the integers  $\{q/2p, \dots, q/2pg\}$  [BGM<sup>+</sup>16].

Note also that for our purposes we require *quasipolynomial*-time hardness (or  $T(n)$ -hardness for  $T(n)$  being any other fixed, mildly superpolynomial function of the security parameter) of LWR. While slightly stronger than standard polynomial-time hardness, this remains a reasonable assumption since algorithms for worst-case lattice problems are still believed to require at least subexponential time.

### 3.2.4 Partial assignments

Let  $\alpha \in \{0, 1, \star\}^d$  be a *partial assignment*. We refer to  $S(\alpha) \subseteq [d] : \alpha_i = \star$  as the set of *free variables* and  $[d] \setminus S(\alpha)$  as the set of *fixed variables*. Given two partial assignments  $\alpha, \beta$ , let the *resolution*  $\alpha \& \beta$  denote the partial assignment  $\gamma$  obtained by substituting  $\alpha$  into  $\beta$ . That is,

$$\gamma_i = \begin{cases} \star & i \in S(\alpha) \setminus S(\beta) \\ \beta_i & i \in [d] \setminus S(\beta) \\ \alpha_i & i \in S(\beta) \setminus S(\alpha) \end{cases} \quad (3.2.3)$$

In this case we say that  $\gamma$  is a *refinement* of  $\beta$  that is the *result of applying*  $\alpha$ . We write  $\gamma \in \text{App}(\alpha)$  to denote that  $\gamma$  is a result of applying  $\alpha$ . Note that the set of refinements of  $\beta$  consists of all  $3^{jS(\beta)}$  partial assignments  $\gamma \in \{0, 1, \star\}^d$  which agree with  $\beta$  on all fixed variables of  $\beta$ .

Given  $\alpha$ , let  $w(\alpha)$  denote  $\sum \alpha_i$ , that is, the Hamming weight of its fixed variables. Note that  $w(\alpha \& \beta) = w(\alpha) + w(\beta)$ .

Given a function  $h : \mathbb{R}^d \rightarrow \mathbb{R}$  and partial assignment  $\gamma$ , we use  $h_\gamma : \mathbb{R}^d \rightarrow \mathbb{R}$  to denote its partial restriction given by substituting in  $\gamma_i$  into the  $i$ -th input coordinate if  $\gamma_i \neq \perp$ . Note that given two partial restrictions  $\alpha, \beta$ ,

$$(h_\beta)_\alpha = h_{\alpha \& \beta} \tag{3.2.4}$$

We say that  $\alpha$  is *sorted* if the restriction of  $\alpha$  to its fixed variables is sorted in nonincreasing order, e.g.  $\alpha = (1, \star, 1, \star, \star, 0, 0)$  is sorted, but  $\alpha = (1, \star, 0, \star, \star, 0, 1)$  is not. Given  $\alpha$  which is not necessarily sorted, denote its *sorting* by  $\bar{\alpha}$ . In general, we will use overline notation to denote sorted partial assignments.

### 3.3 Compressing the Daniely–Vardi Lift

In this section we show how to refine the lifting procedure of Daniely and Vardy [DV21] such that whenever the underlying discrete functions satisfy a property we term *compressibility*, we obtain hardness under the Gaussian for networks with just one extra hidden layer.

**Definition 3.3.1.** Let  $q > 0$  be a modulus.<sup>2</sup> We call an  $L$ -hidden-layer ReLU network  $f : \mathbb{Z}_q^d \rightarrow [0, 1]$  *compressible* if it is expressible in the form  $f(x) = \sigma(h(x))$ , where

- $h : \mathbb{Z}_q^d \rightarrow T$  is an  $(L - 1)$ -hidden-layer network such that  $\|h(x)\| = \text{poly}(d)$  for all  $x$ ;
- $h$  has range  $T = h(\mathbb{Z}_q^d)$  such that  $T \subseteq \mathbb{Z}$  and  $\|T\| = \text{poly}(d)$ ; and
- $\sigma : T \rightarrow [0, 1]$  is a mapping from  $h$ 's possible output values to  $[0, 1]$ .

*Remark 3.3.2.* To see why such an  $f$  is an  $L$ -hidden-layer network in  $z$ , consider the function  $\sigma : T \rightarrow \mathbb{R}$ . Because  $T \subseteq \mathbb{Z}$  and  $\|T\| = \text{poly}(d)$ ,  $\sigma$  is expressible as (the

---

<sup>2</sup>Our results are stronger when  $q$  is taken to be a large polynomial in the dimension, but the Boolean  $q = 2$  case is illustrative of all the main ideas.

restriction to  $T$  of) a piecewise linear function on  $\mathbb{R}$  whose size and maximum slope are  $\text{poly}(d)$ , and hence as a  $\text{poly}(d)$ -sized one-hidden-layer ReLU network from  $\mathbb{R}$  to  $\mathbb{R}$ . By composition,  $x \mapsto \sigma(h(x))$  can be represented by an  $L$ -hidden-layer network.

We now formally state a theorem which captures our “compressed” version of the DV lift. The version of this theorem for  $L + 2$  layers is implicit in [DV21]. In technical terms, our improvement consists of removing the single outer ReLU present in their construction. Thus, while our construction still has three *linear* layers, it has only two *non-linear* layers.

**Theorem 3.3.3** (Compressed DV lift). *Let  $q = \text{poly}(d)$  be a modulus. Let  $C$  be a class of compressible  $L$ -hidden-layer  $\text{poly}(d)$ -sized ReLU networks mapping  $Z_q^d$  to  $[0, 1]$ . Let  $m = m(d) = \omega_d(1)$  be a size parameter that grows slowly with  $d$ . There exists a class  $C^M$  of  $(L + 1)$ -hidden-layer  $d^{(m)}$ -sized ReLU networks mapping  $\mathbb{R}^d$  to  $[0, 1]$  such that the following holds:*

*Suppose there is an efficient algorithm  $A$  capable of learning  $C^M$  over  $N(0, \text{Id}_d)$  up to squared loss  $d^{-(m)}$ . Then there is an efficient algorithm  $B$  capable of weakly predicting  $C$  over  $\text{Unif}(Z_q^d)$  with advantage  $d^{-(m)}$  over guessing the constant  $1/2$  in the following sense: given access to labeled examples  $(x, f(x))$  for  $x \sim \text{Unif}(Z_q^d)$  and an unknown  $f \in C$ ,  $B$  satisfies*

$$\mathbb{E} [(B(x) - f(x))^2] < \mathbb{E} \left[ \left( \frac{1}{2} - f(x) \right)^2 \right] - d^{-(m)},$$

*where the probability is taken over both  $x$  and the internal randomness of  $B$ . We refer to  $C^M$  as the lifted class corresponding to  $C$ .*

By a standard padding argument, we obtain the following corollary which lets us work with polynomial-sized neural networks.

**Corollary 3.3.4** (Compressed DV lift with padding). *Let  $q$ ,  $m$  and  $d$  be as above, and let  $d^0 = d^m$ . View  $C$  and  $C^M$  as function classes on  $Z_q^{d^0}$  and  $\mathbb{R}^{d^0}$  respectively, defined using only the first  $d$  coordinates, so that  $C^M$  is now a  $\text{poly}(d^0)$ -sized class over*

$\mathbb{R}^{d^\theta}$ . Then an algorithm capable of learning  $C^M$  over  $N_{d^\theta}$  up to squared loss  $1/\text{poly}(d^\theta)$  implies a weak predictor for  $C$  over  $\text{Unif}(Z_q^{d^\theta})$  with advantage  $1/\text{poly}(d^\theta)$ .

### 3.3.1 The DV Lift

Before proceeding to the proof of Theorem 3.3.3, we first outline the idea of the original DV lift in the setting of Boolean functions ( $q = 2$ ). The goal is to approximate any given  $f \in C$  by a ReLU network  $f^{\text{DV}} : \mathbb{R}^d \rightarrow \mathbb{R}$  in such a way that  $f^{\text{DV}}$  under  $N_d$  behaves similarly to  $f$  under  $U_d$ . As a first attempt, one might consider the function  $f^*(z) = f(\text{sign}(z))$  (also studied in [KK14a]), where recall that  $\text{sign}(t) = \mathbb{1}[t > 0]$ . We could implement the following reduction: given a random example  $(x, y)$  where  $x \sim U_d$  and  $y = f(x)$ , draw a fresh half-Gaussian  $g \sim \mathcal{N}_d^+$  and output  $((2x - 1)g, y)$  (where the arithmetic in defining the vector  $(2x - 1)g$  is done elementwise). Since  $2x - 1$  is distributed uniformly over  $[-1, 1]$ , the marginal is exactly  $N_d$ , and the labels are consistent with  $f^*$  since  $\text{sign}((2x - 1)g) = x$  and so  $f(\text{sign}((2x - 1)g)) = f(x)$ . However, the issue is that the sign function is discontinuous, and so  $f^*$  is not realizable as a ReLU network.

Daniely and Vardi address this concern by devising a clever construction for  $f^{\text{DV}}$  that interpolates between two desiderata:

- For all but a small fraction of inputs, an initial layer successfully “Booleanizes” the input. In this case, one would like  $f^{\text{DV}}(z)$  to simply behave as  $f(\text{sign}(z))$ .
- For the remaining fraction of inputs, we would ideally like  $f^{\text{DV}}$  to output an uninformative value such as zero, but this would violate continuity of  $f^{\text{DV}}$ .

The trick is to use a continuous approximation of the sign function,  $N_1$ , that interpolates linearly between 0 and 1 on an interval  $[-\delta, \delta]$  (see Fig. 3.1a), and to pair it with a “soft indicator” function  $N_2 : \mathbb{R} \rightarrow \mathbb{R}$  for the region where  $N_1 \neq \text{sign}$ . Concretely,  $N_2(t)$  is constructed as a one-hidden-layer ReLU network that (a) is always nonnegative, (b) equals 0 when  $|t| \geq 2\delta$ , and (c) equals 1 when  $|t| \leq \delta$  (see Fig. 3.1b).



Now let  $N_2^\theta(z) = \sum_j N_2(z_j)$ , and define

$$f^{\text{DV}}(z) = \text{ReLU}(f(N_1(z)) - N_2^\theta(z)). \quad (3.3.1)$$

One can show that  $f^{\text{DV}}$  satisfies  $f^{\text{DV}}(z) = \text{ReLU}(f(\text{sign}(z)) - N_2^\theta(z))$ , since  $N_2^\theta$  “zeroes out”  $f^{\text{DV}}$  wherever  $N_1 \not\in \text{sign}$  for any coordinate. This lets us perform the following reduction: given examples  $(x, y)$  where  $x \sim U_d$  and  $y = f(x)$ , draw a fresh  $g \sim N_d$  and output  $(z, \tilde{y}) = ((2x - 1)g, \text{ReLU}(y - N_2^\theta((2x - 1)g)))$ . The marginal is again  $N_d$ , and the labels are easily seen to be consistent with  $f^{\text{DV}}$ . Correctness of the reduction can be established by using Gaussian anticoncentration to argue that  $f^{\text{DV}}$  is a good approximation of  $f$ . Formally, one can prove the following theorem.

**Theorem 3.3.5** (Original DV lift, implicit in [DV21]). *Let  $\mathcal{C}$  be a class of  $L$ -hidden-layer  $\text{poly}(d)$ -sized ReLU networks mapping  $[0, 1]^d$  to  $[0, 1]$ . There exists a class  $\mathcal{C}^{\text{DV}}$  of  $(L + 2)$ -hidden-layer  $\text{poly}(d)$ -sized ReLU networks mapping  $\mathbb{R}^d$  to  $[0, 1]$  such that the following holds. Suppose there is an efficient algorithm  $A$  capable of learning  $\mathcal{C}^{\text{DV}}$  over  $N(0, \text{Id}_d)$  up to squared loss  $\frac{1}{64}$ . Then there is an efficient algorithm  $B$  capable of weakly predicting  $\mathcal{C}$  over  $\text{Unif}[0, 1]^d$  with squared loss  $\frac{1}{16}$ .*

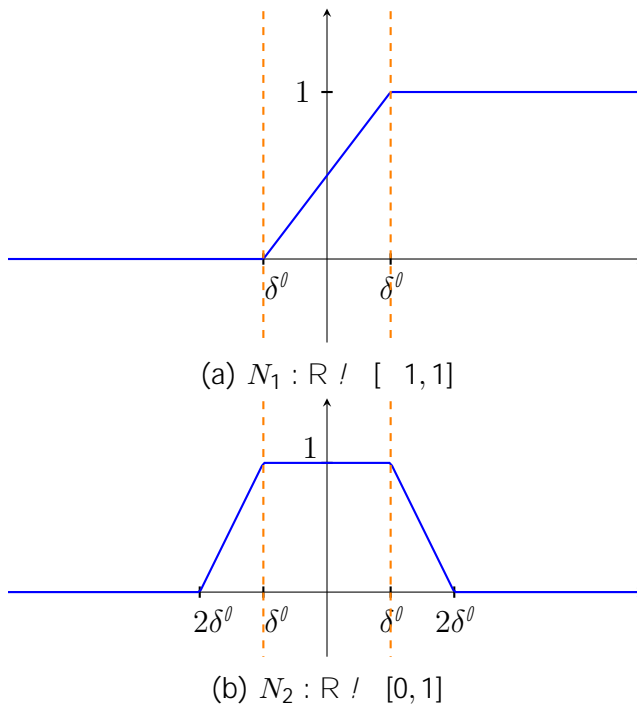


Figure 3.1: Schematic plots of  $N_1$  and  $N_2$  in the  $q = 2$  case, where  $N_2^\theta(z)$  may be realized as  $\sum_{j \in [d]} N_2(z_j)$ . Here,  $\delta^\theta = \Theta(\delta)$  where  $\delta$  is the parameter from Lemmas 3.3.6 and 3.3.7.

We now show how to construct the gadgets  $N_1$  and  $N_2$ , extending them to make them suitable for working with  $Z_q$  for general  $q$  as opposed to just  $\mathcal{N}(0, 1)$ . These constructions utilize the simple but important property that piecewise linear functions on the real line are readily and efficiently realized as linear combination of ReLUs.

Start by letting  $I_0, I_1, \dots, I_{q-1}$  be a partition of  $\mathbb{R}$  into  $q$  consecutive intervals each of mass  $1/q$  under  $\mathcal{N}(0, 1)$  (e.g., when  $q = 2$ ,  $I_0 = (-1, 0)$  and  $I_1 = (0, 1)$ ). Note that these intervals will have differing lengths, and the shortest ones will be the ones closest to the origin. Still, by Gaussian anti-concentration, we know that each  $\int_{I_j} \mathcal{N}(0, 1) \leq \Theta(1/q)$ . Let  $\text{thres}_q : \mathbb{R} \rightarrow Z_q$  be the piecewise constant function that takes on value  $k$  on  $I_k$ . Clearly, when  $t \sim \mathcal{N}(0, 1)$ ,  $\text{thres}_q(t) \sim \text{Unif}(Z_q)$ . Let  $R_1, \dots, R_q$  be intervals such that  $R_k \subset I_{k-1} \cup I_k$  and  $R_k$  contains the boundary point between  $I_{k-1}$

and  $I_k$ , and such that each  $R_k$  has mass  $\delta/q$  for some  $\delta \leq 1$  to be picked later. Let  $S_1, \dots, S_q$  be slightly larger intervals such that  $R_k \subseteq S_k$  for each  $k \in [q]$ , and each  $S_k$  has mass  $2\delta/q$ . By Gaussian anti-concentration again, each  $\mathbb{P}_{z \sim \mathcal{N}(0,1)}[z \in S_k] = \Theta(\delta/q)$ . Notice that by construction,  $\mathbb{P}_{z \sim \mathcal{N}(0,1)}[z \in R_k] = \delta$  and  $\mathbb{P}_{z \sim \mathcal{N}(0,1)}[z \in S_k] = 2\delta$ .

**Lemma 3.3.6.** *Let  $\delta > 0$ ,  $q > 0$ , and intervals  $I_k, R_k, S_k$  for  $k \in [q]$  be as above. There exists a one-hidden-layer ReLU network  $N_1 : \mathbb{R} \rightarrow \mathbb{R}$  with  $O(q)$  units and weights of magnitude  $O(q/\delta)$  such that  $N_1(t) = \text{thres}_q(t)$  if  $t \notin \bigcup_k R_k$ .*

*Proof.* This can be done by considering the piecewise linear function that approximates the function  $\text{thres}_q$  by matching it exactly on  $\mathbb{R} \setminus \bigcup_k R_k$ , and interpolating linearly between values  $k-1$  and  $k$  on the interval  $R_k$  for each  $k \in [q]$ .  $\square$

**Lemma 3.3.7.** *Let  $\delta > 0$ ,  $q > 0$ , and intervals  $I_k, R_k, S_k$  for  $k \in [q]$  be as above. There exists a one-hidden-layer ReLU network  $N_2 : \mathbb{R} \rightarrow [0, 1]$  with  $O(q)$  units and weights of magnitude  $O(q/\delta)$  such that*

$$N_2(t) \text{ is } \begin{cases} = 1 & \text{if } t \in R_k \\ = 0 & \text{if } t \in \mathbb{R} \setminus \bigcup_k S_k \\ 0 & \text{otherwise} \end{cases}.$$

*Proof.* Consider the piecewise linear function that is 0 on  $\mathbb{R} \setminus \bigcup_k S_k$ , is 1 on  $R_k$ , and interpolates linearly between 0 and 1 (or 1 and 0) on  $S_k \setminus R_k$  for every  $k \in [q]$ . Put differently, the graph of  $N_2$  consists of a trapezoid on each  $S_k$  that achieves its maximum value of 1 on  $R_k$ .  $\square$

### 3.3.2 Saving One Hidden Layer via Compressibility

The starting point for exploiting compressibility to avoid a hidden layer in the lift is as follows. Compressibility lets us express  $f(x)$  as  $\sigma(h(x))$  for some  $h : \mathbb{Z}_q^d \rightarrow T$  with a poly( $d$ )-sized range  $T \subseteq \mathbb{Z}$ , and some  $\sigma : T \rightarrow [0, 1]$ . So we can write

$$f(x) = \sigma(h(x)) = \sum_{t \in T} \sigma(t) \mathbb{1}[h(x) = t].$$

We would like a lifted function  $f^M : \mathbb{R}^d \rightarrow \mathbb{R}$  (where we introduce  $f^M$  as notation to distinguish our lift from the original DV lift, denoted  $f^{DV}$ ) such that  $f^M(z)$  behaves like  $\sigma(h(\text{thres}_q(z)))$  except when  $N_2$  indicates that  $N_1 \notin \text{thres}_q$ , in which case we want  $f^M(z) = 0$ . To this end, we might hope to write

$$f^M(z) = \sum_{t \in T} \sigma(t) \mathbb{1}[h(N_1(z)) = t] \mathbb{1}[\exists j : N_2(z_j) < 1].$$

Suppose that we could build a one-hidden-layer network  $N_3(s_1, \dots, s_d; t)$  that behaves like  $\mathbb{1}[t = 0] \mathbb{1}[\exists j : s_j < 1]$ . Then we could realize  $f^M$  as

$$f^M(z) = \sum_{t \in T} \sigma(t) N_3(N_2(z_1), \dots, N_2(z_d); h(N_1(z)) = t).$$

Notice that whenever  $N_2(z_j) = 1$  for any coordinate  $j$ , this expression vanishes. Otherwise, we know that  $h(N_1(z)) = h(\text{thres}_q(z))$ , which takes values in  $T$ , so that only the summand with  $t = h(\text{thres}_q(z))$  survives and the expression simplifies to  $f(\text{thres}_q(z)) N_3(N_2(z_1), \dots, N_2(z_d); 0)$ . It is not hard to show that this is sufficient to let us complete the required reduction. Moreover, because  $N_3$  is a one-hidden-layer network in its arguments, and because both  $h \circ N_1$  and  $N_2$  have at most  $L$  hidden layers (for  $h \circ N_1$ , one comes from  $N_1$  and  $L - 1$  from  $h$ ; for  $N_2$ , it itself has just one hidden layer), this implementation of  $f^M$  would have only  $L + 1$  hidden layers.

Slightly more generally, one can show that it would suffice to build a one-hidden-layer network  $N_3$  with the following properties:

$$N_3(s_1, \dots, s_d; t) = \begin{cases} 0 & \text{if } \exists j : s_j = 1 \\ 0 & \text{if } t \neq 0 \\ 1 & \text{if } \exists j : s_j = 0 \text{ and } t = 0 \end{cases} \quad (3.3.2)$$

Unfortunately, most natural attempts to construct  $N_3$  with such ideal properties — in particular, all formulations of  $N_3$  purely as a function of two variables,  $\sum_j s_j$  and  $t$ , which was the approach taken in [DV21] — run into difficulties and appear to require *two* hidden layers (see Appendix B.1 for discussion). One approach that does almost work, however, comes at the cost of exponential size. Let  $\psi(s_1, \dots, s_d; t)$

be any function that vanishes whenever  $t \geq n$  (for all  $s_1, \dots, s_d \in [0, 1]^d$ ). For simplicity, let us consider the  $d = 3$  case. Consider the following expression that resembles the inclusion-exclusion formula:

$$\psi(s_1, s_2, s_3; t) - \psi(1, s_2, s_3; t) - \psi(s_1, 1, s_3; t) - \psi(s_1, s_2, 1; t) \quad (3.3.3)$$

$$+ \psi(s_1, 1, 1; t) + \psi(1, s_2, 1; t) + \psi(s_1, 1, 1; t) - \psi(1, 1, 1; t) \quad (3.3.4)$$

Notice that whenever any  $s_j = 1$ , this expression vanishes identically. Moreover, for any  $t \geq n$  (and any  $s_1, \dots, s_d$ ), the expression vanishes again because each summand vanishes. Thus the first two properties are satisfied; the third property turns out to be more subtle, and we will ignore it for the moment. The natural generalization of this expression to general  $d$  can be stated in the language of partial assignments.

**Lemma 3.3.8.** *Let  $\psi : \mathbb{R}^d \rightarrow \mathbb{R}$  be any function. Let  $P_i$  denote the set of partial assignments  $\gamma \in \{0, 1\}^d$  with  $i$  1s. The expression*

$$\sum_{i=0}^d \sum_{\gamma \in P_i} (-1)^i \psi_\gamma \quad (3.3.5)$$

*vanishes whenever any  $s_j = 1$ . (We may view  $t$  as an additional parameter that is always left free, as in Eq. (3.3.3))*

*Proof.* For concreteness, suppose  $s_1 = 1$ . Let  $P_i^\star$  (resp.  $P_i^1$ ) denote the set of  $\gamma \in P_i$  with  $s_1 = \star$  (resp.  $s_1 = 1$ ). For every  $i \in \{0, \dots, d-1\}$ , we can form a bijection between  $P_i^\star$  and  $P_{i+1}^1$  using the map  $\gamma \mapsto \gamma^\theta$  where  $\gamma^\theta = (1, \gamma_2, \dots, \gamma_d)$ . When  $s_1 = 1$ , for every such pair  $(\gamma, \gamma^\theta)$ , we have  $\psi_\gamma = \psi_{\gamma^\theta}$ , and moreover they occur in (3.3.5) with opposite signs. Thus the entire expression vanishes.  $\square$

Let us assume for now that  $\psi$  is picked suitably and the rest of the reduction goes through with this construction (as one can verify when we come to the proof of Theorem 3.3.3, this would indeed be the case). This construction has size  $2^d$ ,

meaning that the resulting lifted functions would have size  $S = \text{poly}(2^d)$ . But by Theorem 3.4.5, the SQ lower bound for the LWR functions over  $Z_q^n$  with  $n = d$  and  $q = \text{poly}(n)$  scales as  $q^{(n)} = 2^{(d \log d)} = S^{(\log \log S)}$ , which is still superpolynomial in  $S$ . Thus after padding the dimension to  $d^\theta = 2^d$ , this construction would actually still yield a superpolynomial SQ lower bound for two-hidden-layer ReLU networks over  $\mathbb{R}^{d^\theta}$ .

Instead of pursuing this route, however, we give a more efficient construction that has size only slightly superpolynomial in  $d$ . The key idea is to restrict attention to those possibilities for  $(s_1, \dots, s_d) = (N_2(z_1), \dots, N_2(z_d))$  that are the most likely. Specifically, if  $m = \omega_d(1)$  is the size parameter from Theorem 3.3.3, then by setting  $\delta$  in Lemmas 3.3.6 and 3.3.7 appropriately, we can ensure that with overwhelming probability over  $z \sim N(0, \text{Id})$ , no more than  $m$  of the  $N_2(z_j)$  are simultaneously 1. Accordingly, we focus on constructing  $N_3$  such that

$$N_3(s_1, \dots, s_d; t) = \begin{cases} 0 & \text{if between 1 and } m \text{ of the } s_i \text{ are 1} \\ 0 & \text{if } t \geq \frac{1}{2} \\ 1 & \text{otherwise} \end{cases}. \quad (3.3.6)$$

We now describe a  $d^{(m)}$ -sized construction for  $N_3$  that satisfies the first and second properties exactly, and “approximately” satisfies the third in the sense that it takes on a nonzero value with nonnegligible probability over its inputs. As we will see later, this turns out to be enough for the reduction to go through. The construction retains the spirit of using a linear combination of partial restrictions.

**Lemma 3.3.9** (Main lemma). *Let  $m = m(d) = \omega_d(1)$  be a size parameter. Let  $A$  denote the set of all partial assignments  $\alpha \in \{0, 1, \star\}^d$  for which  $|S(\alpha)| = m$  and  $w(\alpha) = 1$ . Let  $B$  denote the set of all sorted partial assignments given by reordering some element of  $A$  and sorting. Given  $i, j \geq 0$ , let  $B_{i,j}$  denote the set of  $\bar{\beta} \in B$  for which  $|S(\bar{\beta})| = i$  and  $w(\bar{\beta}) = j$ . For any symmetric function  $\psi : \mathbb{R}^d \rightarrow \mathbb{R}$ , define the*

function

$$\psi = \sum_{i=0}^m \sum_{j=1}^{m+1} \binom{m-i}{j} \lambda_{i+j} \sum_{\bar{\beta} \in B_{i,j}} \psi_{\bar{\beta}}, \quad \text{for } \lambda_k = \binom{d-k-1}{m-k+1} \quad (3.3.7)$$

Then

(a)  $\sum_{j=1}^m \binom{d}{m} (d-m) 3^m$

(b)  $\psi$  is symmetric

(c)  $\psi_\alpha : \mathbb{R}^d \rightarrow \mathbb{R}$  is the identically zero function for all  $\alpha \in A$ .

**Lemma 3.3.10.** Let

$$\psi(s_1, \dots, s_d; t) = \sum_{i=1}^d \text{ReLU} \left( t - \left( s_i - \frac{1}{d-1} \sum_{j \neq i} s_j \right) \right) \text{ReLU}(dt),$$

viewed as a function of  $s_1, \dots, s_d$  parameterized by  $t$ , and let  $\psi$  be as above. Define  $N_3(s_1, \dots, s_d; t) = \psi(s_1, \dots, s_d; t)$ . Then

(a)  $N_3(s_1, \dots, s_d; t) = 0$  for any  $t \in \mathbb{R}$  if between 1 and  $m$  of the  $s_j$  are 0

(b)  $N_3(s_1, \dots, s_d; t) = 0$  for any  $s_1, \dots, s_d \in [0, 1]^d$  if  $t \in \mathbb{Z} \setminus \{0\}$

(c)  $N_3$  has size at most  $d^{2m}$

(d)  $N_3(\underbrace{0, \dots, 0}_d, s; 0) = s$  for any  $s \in [0, \frac{1}{d}]$ .

Before proceeding to the proofs of Lemmas 3.3.9 and 3.3.10, let us see how to use them to prove Theorem 3.3.3.

*Proof of Theorem 3.3.3.* For each  $f \in \mathcal{C}$  given by  $f = \sigma \circ h$ , let  $f^M \in \mathcal{C}^M$  be given by

$$f^M(z) = \sum_{t \in T} \sigma(t) N_3(N_2(z_1), \dots, N_2(z_d); h(N_1(z)) - t), \quad (3.3.8)$$

where  $N_1$  and  $N_2$  are from Lemmas 3.3.6 and 3.3.7, with the  $\delta$  parameter set to  $d^{-10m}$ , and  $N_3$  is from Lemma 3.3.10. This is an  $(L + 1)$ -hidden layer network since  $h(N_1)$  and  $N_2$  each have at most  $L$  hidden layers, and  $N_3$  adds an additional layer. By Lemma 3.3.10(c), the size of this network is  $S = d^{-m}$ . Note that whenever  $z$  is such that  $N_2(z_1), \dots, N_2(z_d) < 1$ , then:

- $N_1(z) = \text{thres}_q(z)$ , and so  $h(N_1(z)) = h(\text{thres}_q(z))$  takes only integer values in  $T = h(\mathbb{Z}_q^d)$ ; and
- the only  $t$  for which one of the summands in Eq. (3.3.8) is potentially nonzero is the one given by  $t = h(\text{thres}_q(z))$ .

Thus in this case  $f^M$  simplifies to

$$f^M(z) = \sigma(h(\text{thres}_q(z))) N_3(N_2(z_1), \dots, N_2(z_d); 0) \quad (3.3.9)$$

$$= f(\text{thres}_q(z)) N_3(N_2(z_1), \dots, N_2(z_d); 0). \quad (3.3.10)$$

Further, for  $z$  such that between 1 and  $m$  of the  $N_2(z_j)$  are 1, we know that  $\psi(N_2(z_1), \dots, N_2(z_d); t) = 0$  identically (for all  $t \geq \mathbb{R}$ ), so in this case  $f^M(z) = 0$ . And finally, for  $z$  such that more than  $m$  of the  $N_2(z_j)$  are 1, we have no guarantees on the behavior of  $f^M$ , but as we now show, we have set parameters such that this case occurs only with negligible probability, and we can pretend that 0 is still a valid label in this case. Indeed, by standard Gaussian anti-concentration, for each coordinate  $z_j$  we have  $\mathbb{P}_{z_j}[N_2(z_j) = 1] = \mathbb{P}_{z_j}[z_j \geq [kR_k]] = \delta = d^{-10m}$ . The number of coordinates  $j$  for which  $N_2(z_j) = 1$  thus follows a binomial distribution  $B(d, d^{-10m})$ , which has a decreasing pdf with unique mode at  $b(d+1)d^{-10m} < 0$ . Thus the probability of having at least  $m$  1s is at most

$$\sum_{i=m}^d \binom{d}{i} (d^{-10m})^i (1 - d^{-10m})^{d-i} = (d - m + 1) \binom{d}{m} d^{-10m^2} = dd^m d^{-10m^2} = d^{-9m^2} \quad (3.3.11)$$

for sufficiently large  $d$ . This is negligibly small not only in  $d$  but in the size of the network,  $S = d^{-m}$ .



We now describe the reduction. For each labeled example  $(x, y)$  that the discrete learner  $B$  receives, where  $x \sim \text{Unif}(Z_q^d)$  and  $y = f(x)$  for an unknown  $f \in \mathcal{C}$ ,  $B$  forms a labeled example  $(z, \tilde{y})$  for the Gaussian learner  $A$  as follows. For each coordinate  $j \in [d]$ ,  $z_j$  is drawn from  $N(0, 1)$  conditioned on  $z_j \in I_{x_j}$ . Notice that this way  $\text{thres}_q(z) = x$ , and the marginal distribution on  $z$  is exactly  $N_d$ . The modified label is given by

$$\tilde{y} = \tilde{y}(y, z) = \begin{cases} 0 & \text{if more than } m \text{ of the } N_2(z_j) \text{ are } 1 \\ 0 & \text{if between } 1 \text{ and } m \text{ of the } N_2(z_j) \text{ are } 1 \\ y N_3(N_2(z_1), \dots, N_2(z_d); 0) & \text{otherwise} \end{cases} \quad (3.3.12)$$

Note that in the bottom two cases,  $\tilde{y} = f^M(z)$  exactly; in the top case  $\tilde{y}$  is in general inconsistent with  $f^M$ , but as we have seen, this case occurs with  $\text{negl}(S)$  probability. In particular, with overwhelming probability, no  $\text{poly}(S)$ -time algorithm will ever see non-realizable samples.

So  $B$  can feed these new labeled examples  $(z, \tilde{y})$  to  $A$ . Suppose  $A$  outputs a hypothesis  $\hat{f} : \mathbb{R}^d \rightarrow \mathbb{R}$  such that  $\mathbb{E}_z N_d[(\hat{f}(z) - f^M(z))^2] \leq \epsilon$ . We need to show  $B$  can convert this hypothesis into a nontrivial one for its discrete problem. We first define a “good region”  $G \subseteq \mathbb{R}^d$  where  $f^M$  is guaranteed to be nonzero and nontrivially related to the original  $f$  by saying  $z \in G$  iff  $N_2(z_1), \dots, N_2(z_{d-1}) = 0$ , and  $N_2(z_d) \in (\frac{1}{2d}, \frac{1}{d})$ . Observe that when  $z \in G$ , by Eq. (3.3.10) and Lemma 3.3.10(d) we have

$$f^M(z) = f(\text{thres}_q(z)) N_3(N_2(z_1), \dots, N_2(z_{d-1}), N_2(z_d); 0) \quad (3.3.13)$$

$$= f(x) N_3(0, \dots, 0, N_2(z_d); 0) \quad (3.3.14)$$

$$= y N_2(z_d), \quad (3.3.15)$$

where we use the fact that  $\text{thres}_q(z) = x$ , so that  $f(\text{thres}_q(z)) = f(x) = y$ . Let us compute the probability mass of  $G$ . For coordinates  $j \in [d-1]$ , note that  $\mathbb{P}[N_2(z_j) = 0] = \mathbb{P}[z_j \notin [{}_k S_k]] = 1 - 2\delta = 1 - d^{-m}$ . For  $z_d$ , we need a lower bound on the probability that  $N_2(z_d) \in (\frac{1}{2d}, \frac{1}{d})$ . Consider the behavior of  $N_2$  on just the interval  $S_k$  that is closest to the origin (which will be  $k = dq/2e$ ): it changes linearly from

0 to 1 (and again from 1 to 0) on  $S_k \cap R_k$ . It is not hard to see that  $N_2$  takes values in  $(\frac{1}{2d}, \frac{1}{d})$  on a  $O(1/d)$  fraction of  $S_k$ . Since the Gaussian pdf will be at least some constant on all of  $S_k$ , the probability that  $z_d$  lands in this fraction of  $S_k$  is  $\Omega(jS_k/d) = \Omega(\delta/qd) = d^{-m}$ . Overall, we get that

$$\mathbb{P}[z \geq R] = \mathbb{P}\left[N_2(z_d) \geq \left(\frac{1}{2d}, \frac{1}{d}\right)\right] \prod_{j \in [d-1]} \mathbb{P}[N_2(z_j) = 0] = (1 - d^{-m})^{d-1} d^{-m} = d^{-m},$$

which is still  $1/\text{poly}(S)$  and hence non-negligible in the size  $S$  of the network.

The discrete learner  $B$  can now adapt  $\hat{f}$  as follows. Given a fresh test point  $x \sim \text{Unif}(Z_q^d)$ , the learner forms  $z$  such that for each coordinate  $j \in [d]$ ,  $z_j$  is drawn from  $N(0, 1)$  conditioned on  $z_j \in I_{x_k}$ ; for brevity, we shall denote the random variable  $z$  conditioned on  $x$  (formed in this way) by  $z|x$ . If  $z \geq G$ , then  $B$  predicts  $\hat{y} = \frac{\hat{f}(z)}{N_2(z_d)}$  (recall that when  $z \geq z$ ,  $N_2(z_d) > \frac{1}{2d}$ ), and otherwise it simply predicts  $\tilde{y} = \frac{1}{2}$ . The

square loss of this predictor is given by

$$\mathbb{E}_x \mathbb{E}_{\text{Unif}(Z_d^d)} [(\hat{y} - f(x))^2] = \mathbb{E}_x \mathbb{E}_{z/jx} [(\hat{y} - f(x))^2] \quad (3.3.16)$$

$$= \mathbb{E}_{x,z/jx} [(\hat{y} - f(x))^2 \mathbb{1}_{z \in G}] \mathbb{P}[z \in G] + \mathbb{E}_{x,z/jx} [(\hat{y} - f(x))^2 \mathbb{1}_{z \notin G}] \mathbb{P}[z \notin G] \quad (3.3.17)$$

$$= \mathbb{E}_{x,z/jx} \left[ \left( \frac{\hat{f}(z)}{N_2(z_d)} - f(x) \right)^2 \mathbb{1}_{z \in G} \right] \mathbb{P}[z \in G] + \mathbb{E}_{x,z/jx} \left[ \left( \frac{1}{2} - f(x) \right)^2 \mathbb{1}_{z \notin G} \right] \mathbb{P}[z \notin G] \quad (3.3.18)$$

$$= \mathbb{E}_{x,z/jx} \left[ \left( \frac{\hat{f}(z)}{N_2(z_d)} - \frac{f^M(z)}{N_2(z_d)} \right)^2 \mathbb{1}_{z \in G} \right] \mathbb{P}[z \in G] + \mathbb{E}_x \left[ \left( \frac{1}{2} - f(x) \right)^2 \right] \mathbb{P}[z \notin G] \\ \text{(by Eq. (3.3.15), when } z \in G, f^M(z) = f(x)N_2(z_d)\text{)}$$

$$< 4d^2 \mathbb{E}_z [(\hat{f}(z) - f^M(z))^2 \mathbb{1}_{z \in G}] \mathbb{P}[z \in G] + \mathbb{E}_x \left[ \left( \frac{1}{2} - f(x) \right)^2 \right] \mathbb{P}[z \notin G] \\ \text{(when } z \in G, N_2(z_d) > \frac{1}{2d}\text{)}$$

$$4d^2 \mathbb{E}_z [(\hat{f}(z) - f^M(z))^2] + \mathbb{E}_x \left[ \left( \frac{1}{2} - f(x) \right)^2 \right] \mathbb{P}[z \notin G] \quad (3.3.19)$$

$$4d^2 \epsilon + \mathbb{E}_x \left[ \left( \frac{1}{2} - f(x) \right)^2 \right] \mathbb{P}[z \notin G] \quad (3.3.20)$$

$$= \mathbb{E}_x \left[ \left( \frac{1}{2} - f(x) \right)^2 \right] + 4d^2 \epsilon - \mathbb{E}_x \left[ \left( \frac{1}{2} - f(x) \right)^2 \right] \mathbb{P}[z \in G]. \quad (3.3.21)$$

In the case of the hard classes  $\mathcal{C}$  that we consider, we may assume without loss of generality that  $\mathbb{E}_x \mathbb{E}_{\text{Unif}(Z_d^d)} [(\frac{1}{2} - f(x))^2] = 1/\text{poly}(d)$ , since otherwise the problem of learning  $\mathcal{C}$  is trivial (in fact, in our applications we will have  $\mathbb{E}_x \mathbb{E}_{\text{Unif}(Z_d^d)} [(\frac{1}{2} - f(x))^2] = \Theta(1)$ ). This means that by taking

$$\epsilon = \mathbb{P}[z \in G] / \text{poly}(d) = d^{-\Omega(m)} / \text{poly}(d) = d^{-\Omega(m)}$$

sufficiently small (but still  $1/\text{poly}(S)$ ), we may ensure that the square loss of the discrete learner  $B$  is at most  $\mathbb{E}_x \mathbb{E}_{\text{Unif}(Z_d^d)} [(\frac{1}{2} - f(x))^2] = d^{-\Omega(m)}$ , as desired.  $\square$

*Remark 3.3.11.* The only property of the Gaussian  $N(0, \text{Id}_d)$  used crucially in the proof above is that it is a product distribution  $P = \prod_{i \in [d]} P_i$  where each  $P_i$  is suitably anti-concentrated. By some simple changes to the parameters of  $N_1$ ,  $N_2$  and  $N_3$

(depending on  $P$ ), the proof can be made to work more generally for such distributions  $P$ .

### 3.3.3 Proofs of Lemmas 3.3.9 and 3.3.10

We now detail the proofs involved in the construction of the gadget  $N_3$ .

*Proof of Lemma 3.3.9.* Note that  $jA_j = \binom{d}{m}(d-m)$ . Any partial assignment  $\beta$  has at most  $3^{jS(\beta)}$  refinements, and  $B$  is a subset of all refinements of partial assignments from  $A$ , so  $jB_j = \binom{d}{m}(d-m) 3^m$ .

For the remaining parts of the lemma, it will be useful to observe that  $B$  consists exactly of all partial assignments with  $i$  free variables and  $j$  1s for any  $0 \leq i \leq m$  and  $j \geq 1$  satisfying  $i + j = m + 1$ .

To prove the second part of the lemma, it suffices to show that

$$\sum_{\bar{\beta} \in B_{i,j}} h_{\bar{\beta}} \tag{3.3.22}$$

is symmetric for all  $i, j$ . As transpositions generate the symmetric group on  $d$  elements, it suffices to show that (3.3.22) is invariant under swapping two input coordinates, call them  $a, b \in [d]$ . For all  $\bar{\beta} \in B_{i,j}$  for which  $a, b$  are either both present or both absent in  $S(\bar{\beta})$ , this clearly does not affect the value of  $h_{\bar{\beta}}$ . Now consider the set  $S_a$  (resp.  $S_b$ ) of partial assignments  $\bar{\beta} \in B_{i,j}$  for which only  $a$  (resp. only  $b$ ) is present in  $S(\bar{\beta})$ . There is a clear bijection  $f : S_a \rightarrow S_b$ : given  $\bar{\beta} \in S_a$ , swap the  $a$ - and  $b$ -th entries, and vice-versa, and for any  $\bar{\beta} \in S_a$ , the function  $h_{\bar{\beta}} + h_{f(\bar{\beta})}$  is unaffected by the swapping of input coordinates  $a, b$ . This concludes the proof of the second part of the lemma.

Finally, to prove the third part of the lemma, it suffices to verify it for a single  $\alpha \in A$ , as  $h$  is symmetric. So consider  $\bar{\alpha} = \bar{1}, 0, \dots, 0, \star, \dots, \star, g$ . We apply (3.2.4)

to get

$$h_{\bar{\alpha}} = h_{\bar{\alpha}} \sum_{i=0}^m \sum_{j=1}^{m+1} \binom{m-i}{j} \lambda_{i+j} \sum_{\bar{\beta} \in B_{i,j}} h_{\bar{\alpha} \& \bar{\beta}} \quad (3.3.23)$$

$$= h_{\bar{\alpha}} \sum_{\bar{\gamma} \in B \setminus \text{App}(\alpha)} h_{\bar{\gamma}} \sum_{i=0}^m \sum_{j=1}^{m+1} \binom{m-i}{j} \lambda_{i+j} \sum_{\bar{\beta} \in B_{i,j}} \mathbb{1}[\overline{\bar{\alpha} \& \bar{\beta}} = \bar{\gamma}] \quad (3.3.24)$$

Note that for  $\bar{\gamma} = \bar{\alpha}$ , the only  $\bar{\beta} \in B$  for which  $\overline{\bar{\alpha} \& \bar{\beta}} = \bar{\gamma}$  is  $\bar{\beta} = \bar{\alpha}$ . Indeed, for  $\bar{\beta}$  to be such that  $\overline{\bar{\alpha} \& \bar{\beta}} = \bar{\alpha}$ , it must have  $S(\bar{\beta}) = S(\bar{\alpha})$  and exactly one 1, from which it follows that  $\bar{\beta} = \bar{\alpha}$ . Since  $\bar{\alpha} \in B_{m,1}$ , its coefficient in (3.3.24) is given by

$$\binom{m}{m} \lambda_{m+1} = 1, \quad (3.3.25)$$

and so the  $h_{\bar{\alpha}}$  in (3.3.24) cancels with the  $\bar{\gamma} = \bar{\alpha}$ -th summand in (3.3.24).

In the rest of the proof, we can thus focus on sorted  $\bar{\gamma} \in B \setminus \text{App}(\alpha) \cap \overline{nf\bar{\alpha}g}$ . Note that such  $\bar{\gamma}$  satisfy

$$jS(\bar{\gamma})j < m. \quad (3.3.26)$$

To see this, recall that any  $\bar{\gamma} \in B$  with  $jS(\bar{\gamma})j = m$  must have exactly one 1, and since  $\bar{\gamma} \in \text{App}(\bar{\alpha})$  it must be that  $\bar{\gamma}$  must have  $S(\bar{\gamma}) = S(\bar{\alpha})$  and so  $\bar{\gamma} = \bar{\alpha}$ .

Observe that we must have  $\bar{\gamma}_1 = 1$ . Indeed, it cannot be 0 because  $\bar{\gamma}$  is sorted and has at least one 1. It also cannot be  $\star$ . To see this, consider any  $\bar{\beta}$  for which  $\overline{\bar{\alpha} \& \bar{\beta}} = \bar{\gamma}$ . If we had  $\bar{\beta}_1 \notin \star$ , then clearly  $\bar{\gamma}_1 \notin \star$ . If we had  $\bar{\beta}_1 = \star$ , then  $(\bar{\alpha} \& \bar{\beta})_1 = 1$  (as  $\bar{\alpha}_1 = 1$ ), so  $\bar{\gamma} = \overline{\bar{\alpha} \& \bar{\beta}}$  must also have first entry given by 1.

We are now ready to calculate the coefficient of  $h_{\bar{\gamma}}$  (for each  $\bar{\gamma} \in B \setminus \text{App}(\alpha) \cap \overline{nf\bar{\alpha}g}$ ) in (3.3.24) by adding the coefficients of all the  $\bar{\beta} \in B$  for which

$$\overline{\bar{\alpha} \& \bar{\beta}} = \bar{\gamma}. \quad (3.3.27)$$

First let us consider the contribution of  $\bar{\beta} \in B$  for which  $\bar{\beta}_1 = 1$ . Observe that such  $\bar{\beta}$  must have exactly  $w(\bar{\gamma})$  1s. Furthermore, such a  $\bar{\beta}$  is an element of  $B$  if and

only if it has at most  $m + 1 - w(\bar{\gamma})$  free variables, and the set of free variables in  $\bar{\beta}$  must be  $S(\bar{\gamma}) \sqcup V$  where  $V$  is any subset of  $[d] \setminus \text{flg} \sqcup S(\bar{\alpha})$ . The contribution of all such  $\bar{\beta}$  to the coefficient of  $h_{\bar{\gamma}}$  in (3.3.24) is thus

$$\sum_{i=jS(\bar{\gamma})}^{m+1-w(\bar{\gamma})} \binom{m-i}{1} \lambda_{i+w(\bar{\gamma})} \binom{d-m-1}{i-jS(\bar{\gamma})}, \quad (3.3.28)$$

where here the index  $i$  denotes the total number of free variables in  $\bar{\beta}$ , and the factor of  $\binom{d-m-1}{i-jS(\bar{\gamma})}$  is the number of ways to choose  $V$ .

It remains to consider the contribution from  $\bar{\beta} \geq B$  for which  $\bar{\beta}_1 \neq 1$ . First note that clearly we cannot have  $\bar{\beta}_1 = 0$ , as  $\bar{\beta}$  is sorted and has at least one 1 because it lies in  $B$ . The only possibility is  $\bar{\beta}_1 = \star$ , which we split into two cases based on  $w(\bar{\gamma})$ .

**Case 1:**  $w(\bar{\gamma}) = 1$ . In this case, we claim that there are no  $\bar{\beta} \geq B$  simultaneously satisfying (3.3.27) and  $\bar{\beta}_1 = \star$ . Suppose to the contrary. Then such a  $\bar{\beta}_1$  must have at least one 1 in some other entry (as  $\bar{\beta} \geq B$ ), but this would imply that the resolution  $\bar{\alpha} \& \bar{\beta}$  has at least two 1s, a contradiction. The total coefficient of  $h_{\bar{\gamma}}$  in this case is thus exactly given by (3.3.28). Upon substituting  $w(\bar{\gamma}) = 1$ , this simplifies to

$$\sum_{i=jS(\bar{\gamma})}^{m+1-w(\bar{\gamma})} \binom{m-i}{1} \lambda_{i+1} \binom{d-m-1}{i-jS(\bar{\gamma})} = \sum_{i=jS(\bar{\gamma})}^{m+1-w(\bar{\gamma})} \binom{m-i}{1} \lambda_{i+1} \binom{d-i-2}{d-m-2} \binom{d-m-1}{i-jS(\bar{\gamma})} = 0, \quad (3.3.29)$$

where in the last step we use Lemma B.2.1 (which we can apply because of (3.3.26)).

**Case 2:**  $w(\bar{\gamma}) > 1$ . Observe that we must have  $w(\bar{\beta}) = w(\bar{\gamma}) - 1$  (as the only entry of  $\bar{\alpha}$  equal to 1 is the first entry, and the first entry of  $\bar{\beta}$  is  $\star$ ). As  $w(\bar{\gamma}) - 1 > 0$  in the current case, such a  $\bar{\beta}$  is an element of  $B$  if and only if it has at most  $m + 2 - w(\bar{\gamma})$  free variables, and the set of free variables in  $\bar{\beta}$  must be  $\text{flg} \sqcup S(\bar{\gamma}) \sqcup V$  where  $V$  is any subset of  $[d] \setminus \text{flg} \sqcup S(\bar{\alpha})$ . Thus in this second case, the contribution of all  $\bar{\beta}$

with  $\bar{\beta}_1 = \star$  to the coefficient of  $h_{\bar{\gamma}}$  in (3.3.24) is

$$\sum_{i=jS(\bar{\gamma})j+1}^{m+2-w(\bar{\gamma})} (-1)^{m-i} \lambda_{i+w(\bar{\gamma})-1} \binom{d-m-1}{i \quad jS(\bar{\gamma})j \quad 1} = \sum_{j=jS(\bar{\gamma})j}^{m+1-w(\bar{\gamma})} (-1)^{m-j-1} \lambda_{j+w(\bar{\gamma})} \binom{d-m-1}{j \quad jS(\bar{\gamma})j}, \quad (3.3.30)$$

where here the index  $i$  denotes the total number of free variables in  $\bar{\beta}$ , the factor of  $\binom{d-m-1}{i \quad jS(\bar{\gamma})j \quad 1}$  is the number of ways to choose  $V$  (note that  $jVj = i - jS(\bar{\gamma})j - 1$ ), and in the second expression we made the change of variable  $j = i - 1$ . We conclude that in this case, the coefficient of  $h_{\bar{\gamma}}$  in (3.3.24) is given by the sum of (3.3.28) and (3.3.30), which is 0.

Overall, we conclude that the entire RHS of (3.3.24) vanishes for  $\alpha \geq A$ , proving the third part of the lemma.  $\square$

The next lemma formally constructs  $N_3$  and verifies that it has the required properties, is of acceptable size, and that it takes on nonzero values on a significant part of its domain.

*Proof of Lemma 3.3.10.* Part (a) follows directly from Lemma 3.3.9(c). Part (b) follows by verifying that for any  $t \in \mathbb{Z} \setminus \{0\}$ ,  $\psi(s_1, \dots, s_d; t) = 0$  for any  $s_1, \dots, s_d \in [0, 1]^d$ ; this means that  $\psi$ , which is a combination of partial restrictions of  $\psi$ , also vanishes for such  $t$ . First suppose that  $t$  is a positive integer. Observe that  $t \geq 1$  while  $s_i \leq \frac{1}{d-1} \sum_{j \neq i} s_j \leq [1, 1]$ , so each ReLU in the definition of  $\psi$  is activated and we get

$$\psi(s_1, \dots, s_d; t) = \sum_{i=1}^d \left[ t \cdot \left( s_i - \frac{1}{d-1} \sum_{j \neq i} s_j \right) \right] dt = \sum_{i=1}^d \left( s_i - \frac{1}{d-1} \sum_{j \neq i} s_j \right) = 0. \quad (3.3.31)$$

Next suppose that  $t$  is a negative integer. Then  $t \leq -1$  while  $s_i \leq \frac{1}{d-1} \sum_{j \neq i} s_j \leq [1, 1]$ , so each ReLU in the definition of  $h$  is inactive and we get  $\psi(s_1, \dots, s_d; t) = 0$ .

For part (c), observe that by the size bound in Lemma 3.3.9(a) and the fact that  $\psi$  contains  $O(d)$  ReLUs, the size of  $N_3$  may be bounded by

$$S = O(d) \binom{d}{m} (d - m) (3^m + 1) = O(d) \left( \frac{d^{m+1} 3^m}{m!} + 1 \right) = d^{m+2} = d^{2m}$$

for  $m$  larger than some absolute constant.

It remains to prove part (d). For brevity, we will omit the parameter  $t$  and just refer to  $\psi(0, \dots, 0, s; t)$  and  $\psi(0, \dots, 0, s; t)$  as  $\psi(0, \dots, 0, s)$  and  $\psi(0, \dots, 0, s)$ . We first compute  $\psi(0, \dots, 0, s)$ : for  $s \geq [0, 1]$ ,

$$\psi(0, \dots, 0, s) = \text{ReLU}(s) + (d - 1) \text{ReLU}\left(\frac{1}{d - 1} s\right) = s. \quad (3.3.32)$$

Next, for any  $\bar{\beta} \geq B$ , if  $w(\bar{\beta}) = j$  for some  $0 \leq j \leq m + 1$ , then if  $\bar{\beta}_d = \star$ ,

$$\psi_{\bar{\beta}}(0, \dots, 0, s) \quad (3.3.33)$$

$$= \psi(\underbrace{1, \dots, 1}_j, \underbrace{0, \dots, 0}_{d-j-1}, s) \quad (3.3.34)$$

$$= j \text{ReLU}\left(1 + \frac{1}{d - 1}(j - 1 + s)\right) + (d - j - 1) \text{ReLU}\left(\frac{1}{d - 1} j + \frac{1}{d - 1} s\right) \quad (3.3.35)$$

$$+ \text{ReLU}\left(s + \frac{1}{d - 1} j\right) \quad (3.3.36)$$

$$= \frac{d - j - 1}{d - 1} (j + s) + \text{ReLU}\left(s + \frac{1}{d - 1} j\right) \quad (3.3.37)$$

Note that when  $s \geq [0, 1/(d - 1)]$ , because  $j - 1$  (as  $\bar{\beta} \geq B$ ) this simplifies to

$$= \frac{(d - j - s)j}{d - 1}. \quad (3.3.38)$$

On the other hand, if  $\bar{\beta}_d \geq \bar{f}0, 1g$ , then

$$\psi_{\bar{\beta}}(0, \dots, 0, s) = \psi(\underbrace{1, \dots, 1}_j, \underbrace{0, \dots, 0}_{d-j}, s) \quad (3.3.39)$$

$$= j \text{ReLU}\left(1 + \frac{1}{d - 1}(j - 1)\right) + (d - j) \text{ReLU}\left(\frac{1}{d - 1} j\right) = \frac{(d - j)j}{d - 1}. \quad (3.3.40)$$



As there are  $\binom{d-1}{i}$  (resp.  $\binom{d-1}{i}$ ) partial assignments in  $B_{i,j}$  for which  $\bar{\beta}_d = \star$  (resp.  $\bar{\beta}_d \geq f(0,1g)$ ), we can thus explicitly compute  $h(0, \dots, 0, s)$  for  $s \in [0, 1/(d-1)]$  to be

$$\psi(0, \dots, 0, s) = \sum_{i=0}^m \sum_{j=1}^{m+1-i} \binom{m}{i} \binom{d-i-j-1}{m-i-j+1} \left( \binom{d-1}{i-1} \frac{(d-j-s)j}{d-1} + \binom{d-1}{i} \frac{(d-j)j}{d-1} \right). \quad (3.3.41)$$

By Lemma B.2.2, the double sum is equal to zero, so  $h(0, \dots, 0, s) = h(0, \dots, 0, s) = s$  for  $s \in [0, 1/(d-1)]$  as claimed.  $\square$

### 3.4 Statistical Query Lower Bound

We prove a superpolynomial SQ lower bound (for general queries as opposed to only correlational or Lipschitz queries) for weakly learning two-hidden-layer ReLU networks under the standard Gaussian.

**Theorem 3.4.1.** *Fix any  $\alpha \in (0, 1)$ . Any SQ learner capable of learning  $\text{poly}(d)$ -sized two-hidden-layer ReLU networks under  $N(0, \text{Id}_d)$  up to squared loss  $\epsilon$  (for some sufficiently small  $\epsilon = 1/\text{poly}(d)$ ) using bounded queries of tolerance  $\tau \geq 2^{(\log d)^2 - \alpha}$  must use at least  $\Omega(2^{2^{(\log d)^\alpha}} \tau^2) = d^{\omega(1)} \tau^2$  such queries.*

For instance, taking  $\alpha = \frac{1}{2}$  gives a slightly subexponential (but super-quasipolynomial) in  $d$  query lower bound for queries of tolerance at least inverse quasipolynomial in  $d$ .

This theorem is proven using the following key reduction, which adapts the compressed DV lift (Theorem 3.3.3) to the SQ setting.

**Theorem 3.4.2.** *Let  $q = \text{poly}(d)$  be a modulus, and let  $m = m(d) = \omega_d(1)$  be a size parameter. Let  $\mathcal{C}$  be a class of compressible  $L$ -hidden-layer  $\text{poly}(d)$ -sized ReLU networks mapping  $\mathbb{Z}_q^d$  to  $[0, 1]$ , and let  $\mathcal{C}^M$  be the lifted class of  $(L+1)$ -hidden-layer  $d^{(m)}$ -sized ReLU networks corresponding to  $\mathcal{C}$ , mapping  $\mathbb{R}^d$  to  $\mathbb{R}$  (as in Theorem 3.3.3). Suppose there is an SQ learner  $A$  capable of learning  $\mathcal{C}^M$  over  $N(0, \text{Id}_d)$  up to squared loss  $d^{-(m)}$  using queries of tolerance  $\tau$ , where  $\tau \geq d^{-(m^2)}$ . Then there is an SQ learner  $B$  that, using the same number of queries of tolerance  $\tau/2$ , produces a weak*

predictor  $\tilde{B}$  for  $C$  over  $\text{Unif}(Z_q^d)$  with advantage  $d^{-\Omega(m)}$  over guessing the constant  $1/2$  (in expectation over both the data and the internal randomness of  $\tilde{B}$ ).

*Proof.* Recall that  $B$  is given SQ access to a distribution of pairs  $(x, y)$  where  $x \in \text{Unif}(Z_q^d)$  and  $y = f(x)$  for an unknown  $f \in C$ .  $A$  can request estimates  $\mathbb{E}[\phi(x, y)] \pm \tau$  for arbitrary bounded queries  $\phi : Z_q^d \times [0, 1] \rightarrow [0, 1]$  and any desired  $\tau$ . We know that given  $(x, y)$ , the distribution of  $(z, \tilde{y})$ , where  $z = z(x)$  is defined by drawing each  $z_j$  from  $\mathcal{N}(0, 1)$  conditioned on  $z_j \geq I_{x_j}$  and  $\tilde{y} = \tilde{y}(y, z)$  is as in Eq. (3.3.12), is consistent with some  $f^M \in C^M$  except on a region of probability mass at most  $d^{-\Omega(m^2)}$  (recall Eq. (3.3.11)). Suppose we could simulate SQ access to the distribution of  $(z, f^M(z))$  using only SQ access to that of  $(x, f(x))$ . Then by the argument in Theorem 3.3.3, simulating  $A$  on the  $(z, f^M(z))$  distribution would give us a weak predictor  $\tilde{B}$  for the distribution of  $(x, f(x))$ , satisfying

$$\mathbb{E}[(\tilde{B}(x) - f(x))^2] < \mathbb{E}[(\frac{1}{2} - f(x))^2] + d^{-\Omega(m)}.$$

What we must describe is how  $B$  can simulate  $A$ 's statistical queries. Say  $A$  requests an estimate  $\mathbb{E}_z[\phi(z, f^M(z))] \pm \tau$  for some query  $\phi : \mathbb{R}^d \times \mathbb{R} \rightarrow [0, 1]$ . Consider the query  $\tilde{\phi} : Z_q^d \times [0, 1] \rightarrow [0, 1]$  given by  $\tilde{\phi}(x, y) = \mathbb{E}_{z(x)}[\phi(z(x), \tilde{y}(y, z(x)))]$ . This function can be computed without any additional SQs, since the distribution of  $(z, \tilde{y}) = (z(x), \tilde{y}(y, z(x)))$ , given  $(x, y)$ , is fully determined and known to  $B$ . Observe that

$$\mathbb{E}_{x,y}[\tilde{\phi}(x, y)] = \mathbb{E}_{x,z(x)}[\phi(z(x), \tilde{y}(y, z(x)))] = \mathbb{E}_{z,\tilde{y}}[\phi(z, \tilde{y})]. \quad (3.4.1)$$

We must also account for the difference between  $\mathbb{E}_z[\phi(z, f^M(z))]$  and  $\mathbb{E}_{z,\tilde{y}}[\phi(z, \tilde{y})]$ . But because the distributions only differ on a region of mass  $d^{-\Omega(m^2)}$  and  $\phi$  is bounded, we have

$$\left| \mathbb{E}_z[\phi(z, f^M(z))] - \mathbb{E}_{z,\tilde{y}}[\phi(z, \tilde{y})] \right| = \Theta(d^{-\Omega(m^2)}) \leq \frac{\tau}{2} \quad (3.4.2)$$

since we assumed  $\tau = d^{-\Omega(m^2)}$ . Putting together (3.4.1) and (3.4.2), we see that  $B$  can simulate  $A$ 's query  $\phi$  to within tolerance  $\tau$  by querying  $\tilde{\phi}$  with tolerance  $\tau/2$ .  $\square$

Again, by a padding argument we can obtain a corollary similar to Corollary 3.3.4, for which we omit the formal statement. We will use such an argument in the proof of Theorem 3.4.1.

### 3.4.1 SQ lower bound via parities

We can obtain an SQ lower bound for two-hidden-layer ReLU networks by lifting the problem of learning parities under  $U_d$ , which is well-known to require exponentially many queries. More precisely, we show that an SQ learner for two-hidden-layer ReLU networks would yield an SQ algorithm for the problem of distinguishing an unknown parity from random labels.

**Theorem 3.4.3** ([Kee98, BFJ<sup>+</sup>94]). *Consider an SQ algorithm given SQ access either to the distribution of labeled pairs  $(x, y)$  where  $x \sim U_d$  and  $y = \chi_S(x)$  for an unknown parity  $\chi_S$  or to the randomly labeled distribution  $U_d \sim \text{Unif}\{0, 1\}^d$ . Any algorithm capable of distinguishing between the two cases with probability  $2/3$  using queries of tolerance  $\tau$  requires at least  $\Omega(2^d \tau^2)$  such queries.*

**Lemma 3.4.4.** *For every  $S \subseteq [d]$ , the parity function  $\chi_S : \{0, 1\}^d \rightarrow \{0, 1\}$  can be implemented as a compressible one-hidden-layer ReLU network of  $\text{poly}(d)$  size.*

*Proof.* Recall that  $\chi_S(x)$  evaluates to 1 if the Hamming weight of the bits of  $x$  in  $S$  is odd, and 0 otherwise, so that  $\chi_S(x) = \sigma(\sum_{j \in S} x_j)$ . This satisfies the definition of a compressible one-hidden-layer network with the inner depth-0 network being  $x \mapsto \sum_{j \in S} x_j$  and  $\sigma(t) = \mathbb{1}[t \text{ is odd}]$ .  $\square$

We can now supply one proof of Theorem 3.4.1.

*First proof of Theorem 3.4.1.* Let  $m = m(d) = \log^c d$  for  $c = \frac{1}{\alpha} - 1$ , and let  $d^\flat = d^m = 2^{\log^{c+1} d}$ , so that  $d = 2^{\log^{1/(1+c)} d^\flat}$ . By Lemma 3.4.4, the class  $\mathcal{C}$  of parities on  $\{0, 1\}^d$  can be implemented by compressible one-hidden-layer  $\text{poly}(d)$ -sized ReLU networks, and so the lifted class  $\mathcal{C}^M$  can be implemented by two-hidden-layer  $d^{(m)}$ -sized ReLU

networks over  $\mathbb{R}^d$ . A padding argument lets us embed these classes into dimension  $d^\theta$ . By using the predictor from Theorem 3.4.2 (with  $q = 2$ ), we could obtain an SQ algorithm capable of distinguishing parities from random labels using queries of tolerance  $\tau/2$ , assuming  $\tau \leq d^{-\alpha}$  with  $\alpha = \frac{1}{1+c}$ . By Theorem 3.4.3, the lower bound for learning parities is  $\Omega(2^d \tau^2) = \Omega(2^{2^{\log^{1/(1+c)} d^\theta}} \tau^2)$ . Substituting  $\alpha = \frac{1}{1+c}$  gives the result.  $\square$

But the SQ lower bound obtained this way via parities is somewhat unconvincing since there is a non-SQ algorithm capable of learning the lifted function class obtained from parities. Indeed, suppose we are given examples  $(z, f^M(z))$  where  $f$  is an unknown parity. We know that whenever  $z$  lands in the “good region”  $G$  from the proof of Theorem 3.3.3 (which happens with non-negligible probability), we have  $f^M(z) = f(\text{sign}(z))N_2(z)$  (recall Eq. (3.3.15)). This means we can simply filter out all  $z \notin G$  and form a clean data set of labeled points  $(\text{sign}(z), f(\text{sign}(z)))$ . The unknown  $f$  (and hence  $f^M$ ) can now be learnt by simple Gaussian elimination. In order to give a more convincing lower bound, we now provide an alternative proof based on LWR.

### 3.4.2 SQ lower bound via the LWR functions

Here we provide an alternative proof of Theorem 3.4.1 using the LWR functions. The hard function class obtained this way is not only *unconditionally* hard for SQ algorithms, it is arguably hard for non-SQ algorithms as well, since LWR is believed to be cryptographically hard.

We begin by stating an SQ lower bound for the LWR functions. This theorem is proven in Section B.3 using a general formulation in terms of pairwise independent function families that may be of independent interest, communicated to us by Bogdanov [Bog21].

**Theorem 3.4.5.** *Let  $C_{LWR}$  denote the  $LWR_{n,p,q}$  function class. Any SQ learner capable of learning  $C_{LWR}$  up to squared loss  $1/16$  under  $\text{Unif}(Z_q^n)$  using queries of tolerance  $\tau$  requires at least  $\Omega(q^n \tau^{-2})$  such queries.*

The following lemma shows that the LWR functions may be realized as compressible one-hidden-layer ReLU networks.

**Lemma 3.4.6.** *For every  $w \in \mathbb{Z}_q^n$ , the LWR function  $f_w : \mathbb{Z}_q^n \rightarrow \mathbb{Z}_p/p$  can be implemented as a compressible one-hidden-layer ReLU network of size  $O(q^2n)$ .*

*Proof.* By definition, we have  $f_w(x) = \frac{1}{p}b(w \cdot x) \bmod qe_p$ , which is a compressible one-hidden-layer ReLU network with the inner depth-0 network (i.e., affine function) being  $w \cdot x$  and  $\sigma(t) = \frac{1}{p}bt \bmod qe_p$ . The size bound follows by observing that for any  $x \in \mathbb{Z}_q^n$ , the quantity  $w \cdot x$  is an integer in  $\{0, \dots, q^2ng\}$ .  $\square$

We are ready for an alternative proof of Theorem 3.4.1.

*Alternative proof of Theorem 3.4.1.* Let  $n$  be the security parameter, and fix moduli  $p, q \geq 1$  such that  $p, q = \text{poly}(n)$  and  $p/q = \text{poly}(n)$ . Let  $d = n$ , so that the SQ lower bound from Theorem 3.4.5 is  $\Omega(q^{n-1}) = d^{-\alpha} = 2^{-\tilde{c}d}$ . Let  $m = m(d) = \log^c d$  for  $c = \frac{1}{\alpha} - 1$ , and let  $d^\theta = d^m = 2^{\log^{c+1} d}$ , so that  $d = 2^{\log^{1/(1+c)} d^\theta}$ . By Lemma 3.4.6, the  $\text{LWR}_{n,p,q}$  function class  $\mathcal{C}_{\text{LWR}}$  is implementable by one-hidden-layer ReLU networks over  $\mathbb{Z}_q^d$  of size  $\text{poly}(n) = \text{poly}(d)$ . The result now follows by Theorem 3.4.2 and the same padding argument as in the proof based on parities.  $\square$

### 3.5 Cryptographic Hardness Based on LWR

In this section we show hardness of learning two-hidden-layer ReLU networks over Gaussian inputs based on LWR. This is a direct application of the compressed DV lift (Theorem 3.3.3) to the LWR problem, which is by definition a hard learning problem over  $\text{Unif}(\mathbb{Z}_q^n)$ , or equivalently  $\text{Unif}(\mathbb{Z}_q^d)$  with  $d = n$ .

**Theorem 3.5.1.** *Let  $n$  be the security parameter, and  $x$  moduli  $p, q \geq 1$  such that  $p, q = \text{poly}(n)$  and  $p/q = \text{poly}(n)$ . Let  $d = n$ . Let  $c > 0$ ,  $m = m(d) = \log^c d$  and*

$d^\theta = d^m$ . Suppose there exists a  $\text{poly}(d^\theta)$ -time algorithm capable of learning  $\text{poly}(d^\theta)$ -sized depth-2 ReLU networks under  $N(0, \text{Id}_{d^\theta})$  up to squared loss  $1/\text{poly}(d^\theta)$ . Then there exists a  $\text{poly}(d^\theta) = 2^{(\log^{1+c} n)}$  time algorithm for  $\text{LWR}_{n,p,q}$ .

*Proof of Theorem 3.5.1.* By Lemma 3.4.6, we know that the class  $\mathcal{C}_{\text{LWR}}$  is implementable by compressible  $\text{poly}(d)$ -sized one-hidden-layer ReLU networks over  $Z_q^d$ , or, after padding, over  $Z_q^{d^\theta}$ . Let  $\mathcal{C}_{\text{LWR}}^M$  denote the corresponding lifted class of  $\text{poly}(d^\theta)$ -sized two-hidden-layer ReLU networks, padded to have domain  $\mathbb{R}^{d^\theta}$ . Applying Corollary 3.3.4 to the assumed learner for  $\mathcal{C}_{\text{LWR}}^M$ , we obtain a  $\text{poly}(d^\theta)$ -time weak predictor predictor for  $\mathcal{C}_{\text{LWR}}$ , which readily yields a corresponding distinguisher for the  $\text{LWR}_{n,p,q}$  problem. Using the facts that  $d^\theta = d^m = 2^{\log^{1+c} d}$  and  $d = n$ , we may translate  $\text{poly}(d^\theta)$  into  $2^{(\log^{1+c} n)}$ , yielding the result.  $\square$

*Remark 3.5.2.* Note that the choice of  $m = m(d) = \log^c d$  in Theorem 3.5.1 is purely for simplicity. By picking  $m(d) = \omega_d(1)$  to be a suitably slow-growing function of  $d$ , such as  $\log d$ , we can obtain a running time for the final LWR algorithm that is as mildly superpolynomial as we like.

In addition, as an immediate corollary of Lemma 3.4.6, we also obtain a hardness result for one-hidden-layer networks under  $\text{Unif } \mathcal{F}_0, 1\mathcal{G}^d$ , improving on the hardness result of [DV21] (see Theorem 3.4 therein) for two-hidden-layer networks under  $\text{Unif } \mathcal{F}_0, 1\mathcal{G}^d$ . For this application, we let  $d = n \log q = \tilde{O}(n)$ , so that we may identify the domain  $Z_q^n$  with  $\mathcal{F}_0, 1\mathcal{G}^d$  via the binary representation. This also identifies  $\text{Unif}(Z_q^n)$  with  $\text{Unif } \mathcal{F}_0, 1\mathcal{G}^d$ .

**Corollary 3.5.3.** *Let  $n, p, q$  be such that  $p, q = \text{poly}(n)$  and  $p/q = \text{poly}(n)$ , and let  $d = n \log q = \tilde{O}(n)$ . Suppose there exists an efficient algorithm for learning  $\text{poly}(d)$ -sized one-hidden-layer ReLU networks under  $U_d$  up to squared loss  $1/4$ . Then there exists an efficient algorithm for  $\text{LWR}_{n,p,q}$ .*

### 3.6 Hardness of Learning using Label Queries

The main result of this section is to show hardness of learning constant-depth ReLU networks over Gaussians from label queries:

**Theorem 3.6.1.** *Assume there exists a family of PRFs mapping  $\{0, 1\}^d$  to  $\{0, 1\}^g$  implemented by  $\text{poly}(d)$ -sized  $L$ -hidden-layer ReLU networks. Then there does not exist an efficient learner that, given query access to an unknown  $\text{poly}(d)$ -sized  $(L+2)$ -hidden-layer ReLU network  $f : \mathbb{R}^d \rightarrow \mathbb{R}$ , is able to output a hypothesis  $h : \mathbb{R}^d \rightarrow \mathbb{R}$  such that  $\mathbb{E}_{z \sim N(0, \text{Id}_d)}[(h(z) - f(z))^2] \leq 1/16$ .*

We first recall the classical connection between pseudorandom functions and learning from label queries (also known as membership queries in the Boolean setting), due to Valiant [Val84] (see e.g. [BR17, Proposition 12] for a modern exposition).

**Lemma 3.6.2.** *Let  $\mathcal{C} = \{f_s\}_{s \in \mathcal{S}}$  be a family of PRFs from  $\{0, 1\}^d$  to  $\{0, 1\}^g$  indexed by the key  $s$ . Then there cannot exist an efficient learner  $L$  that, given query access to an unknown  $f_s \in \mathcal{C}$ , satisfies*

$$\mathbb{P}_{x,s}[L(x) = f_s(x)] \leq \frac{1}{2} + \frac{1}{\text{poly}(d)},$$

where the probability is taken over the random key  $s$ , the internal randomness of  $L$ , and a random test point  $x \sim \text{Unif}\{0, 1\}^d$ .

There exist multiple candidate constructions of PRF families in the class  $\text{TC}^0$  of constant-depth Boolean circuits built with AND, OR, NOT and threshold (or equivalently majority) gates. Because the majority gate can be simulated by a linear combination of ReLUs similar to  $N_1$  from Lemma 3.3.6, any  $\text{TC}_L^0$  (meaning depth- $L$ ) function  $f : \{0, 1\}^d \rightarrow \{0, 1\}^g$  may be implemented as a  $\text{poly}(d)$ -sized  $L$ -hidden-layer ReLU network (see e.g. [VRPS21, Lemma A.3]<sup>3</sup>). Thus we may leverage the following candidate PRF constructions in  $\text{TC}^0$  for our hardness result:

---

<sup>3</sup>Note that what the authors term a depth- $(L+1)$  network is in fact an  $L$ -hidden-layer network in our terminology.

- PRFs in  $\text{TC}_4^0$  based on the decisional Diffie-Hellman (DDH) assumption [KL01] (improving on [NR97]), yielding hardness for depth-6 ReLU networks
- PRFs in  $\text{TC}^0$  based on Learning with Errors [BPR12, BP14], yielding hardness for depth- $O(1)$  ReLU networks

Note that depth 4 is the shallowest depth for which we have candidate PRF constructions based on widely-believed assumptions, and the question of whether there exist PRFs in  $\text{TC}_3^0$  is a longstanding open question in circuit complexity [Raz92, HMP<sup>+</sup>93, RR97, KL01]. Under less widely-believed assumptions, [BIP<sup>+</sup>18] have also proposed candidate PRFs in  $\text{ACC}_3^0$ .

We can now complete the proof of Theorem 3.6.1. Since pseudorandom functions are not necessarily compressible, we will simply use the original DV lift (Theorem 3.3.5).

*Proof of Theorem 3.6.1.* Let  $f_s : \{0,1\}^d \rightarrow \{0,1\}^d$  be an unknown  $L$ -hidden-layer ReLU network obtained from the PRF family by picking the key  $s$  at random. Consider the lifted  $(L+2)$ -hidden-layer ReLU network  $f_s^{\text{DV}} : \mathbb{R}^d \rightarrow \mathbb{R}^d$  from Eq. (3.3.1), given by  $f_s^{\text{DV}}(z) = \text{ReLU}(f_s(N_1(z)) \oplus N_2^\theta(z))$ , where  $N_1$  and  $N_2$  are from Lemmas 3.3.6 and 3.3.7, and  $N_2^\theta(z) = \sum_j N_2(z_j)$ . Suppose there were an efficient learner  $A$  capable of learning functions of the form  $f_s^{\text{DV}}$  using queries. By the DV lift (Theorem 3.3.5),  $A$  yields an efficient predictor  $B$  achieving small constant error w.r.t. the unknown  $f_s$ , contradicting Lemma 3.6.2. We only need to verify that  $A$ 's query access to  $f_s^{\text{DV}}$  can be simulated by  $B$ . Indeed, suppose  $A$  makes a query to  $f_s^{\text{DV}}$  at a point  $z \in \mathbb{R}^d$ . Then  $B$  can make a query to  $f_s$  at the point  $\text{sign}(z)$  and return  $\text{ReLU}(f_s(\text{sign}(z)) \oplus N_2^\theta(z)) = f_s^{\text{DV}}(z)$ , as this was the key property satisfied by  $f_s^{\text{DV}}$ . This completes the reduction and proves the theorem.  $\square$



# Chapter 4: Statistical-Query Lower Bounds via Functional Gradients

## 4.1 Introduction

In this chapter 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<sup>+</sup>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  $H$  be a function class on  $\mathbb{R}^n$ , and let  $D$  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  $H$  under  $D$  with error  $\epsilon$  if it outputs a function  $f$  such that

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

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 Sections C.5 and C.6 for details.)

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

**Theorem 4.1.1.** *Let  $H_{\text{ReLU}}$  be the class of ReLUs on  $\mathbb{R}^n$  with unit weight vectors. Suppose that there is an SQ learner capable of learning  $H_{\text{ReLU}}$  under  $D$  with error  $\epsilon$  using  $q(n, \epsilon, \tau)$  queries of tolerance  $\tau$ . Then for any  $\epsilon$ , 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 4.1.2.** *Consider the above setup with  $H_\sigma$ , the class of unit-weight sigmoid units on  $\mathbb{R}^n$ . For any  $\epsilon$ , there exists  $\tau = n^{-(\log^2 1/\epsilon)}$  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  $n^{O(1/\epsilon^2)}$  [DGJ<sup>+</sup>10]. The only known hardness result for this problem is due to Klivans and Kothari [KK14a] 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 4.1.3.** *Consider the above setup with  $H_{\text{hs}}$ , the class of unit-weight halfspaces on  $\mathbb{R}^n$ . For any  $\epsilon$ , there exists  $\tau = n^{-(1/\epsilon)}$  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  $\tau$ , our constraint on  $\tau$  here can be interpreted as saying that to avoid the exponential query

lower bound, one needs sample complexity at least  $\Theta(1/\tau^2) = n^{\Omega(1/\epsilon)}$ , nearly matching the upper bound of [KKMS08, DGJ<sup>+</sup>10].

These results are formally stated and proved in Section 4.5. More generally, we show in Section 4.6 that our results give superpolynomial SQ lower bounds for agnostically learning any non-polynomial activation. (See Section C.1 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, VW19a], 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 4.7 that in the agnostic setting, no efficient SQ algorithm exists.

**Theorem 4.1.4.** *Consider the above setup with  $H_{\text{mon}}$ , the class of multilinear monomials of degree at most  $d$  on  $\mathbb{R}^n$ . For any  $\epsilon = \exp(-\Theta(d))$  and  $\tau = \epsilon^2$ ,  $q(n, \epsilon, \tau) = n^{\Omega(d)} \tau^{5/2}$ .*

### 4.1.1 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

$\phi$  can be used as a subroutine for learning depth-two neural networks of the form  $\psi(\sum_i \phi(w^i \cdot x))$  where  $\psi$  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 [KS94b]. 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 [KK14a]).

#### 4.1.2 Related Work

Several recent papers have considered the computational complexity of learning simple neural networks [Bac17, GKKT17a, YS20, FCG20, KK14a, LSSS14, SVWX17, VW19a, 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. [DKZ20b] 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  $\epsilon$  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  $\epsilon$  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.

### 4.1.3 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 4.3 we explain our usage of functional gradient descent, with Assumption 4.3.1 formally stating the kind of agnostic learning guarantee we require for a single neuron. The main reduction itself is Theorem 4.4.1, the subject of Section 4.4. In Sections 4.5, 4.6 and 4.7 we derive the formal lower bounds which follow as a consequence of our reduction. Finally in Section 4.8, we contrast these lower bounds by also including some simple upper bounds.

## 4.2 Preliminaries

**Notation** Let  $D$  be a distribution over  $\mathbb{R}^n$ , which for us will be the standard Gaussian  $N(0, I_n)$  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  $B_D(R) = \{f \in L^2(\mathbb{R}^n, D) : \|f\|_D \leq R\}$ . 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  $\{0, 1\}$ -valued label  $y$  in such a way that  $\mathbb{E}[y|x] = f(x)$  for a fixed function  $f : \mathbb{R}^n \rightarrow [0, 1]$ , known as its conditional mean function. We will use  $D_f$  to refer to the (unique) induced labeled distribution on  $\mathbb{R}^n \times \{0, 1\}$ , 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  $y \in \{0, 1\}$  and  $\mathbb{E}[y|x] = f(x)$ .

#### 4.2.1 Statistical Query (SQ) Model

A statistical query is specified by a query function  $\phi : \mathbb{R}^n \times \{0, 1\} \rightarrow [0, 1]$ . Given a labeled distribution  $D$  on  $\mathbb{R}^n \times \{0, 1\}$ , the SQ model allows access to an SQ oracle (known as the STAT oracle in the SQ literature) that accepts a query  $\phi$  of specified tolerance  $\tau$ , and responds with a value in  $[\mathbb{E}_{(x,y) \sim D}[\phi(x, y)] - \tau, \mathbb{E}_{(x,y) \sim D}[\phi(x, y)] + \tau]$ . One can interpret the tolerance  $\tau$  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  $\tau$ , 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  $\epsilon$  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  $\hat{f} : \mathbb{R}^n \rightarrow [0, 1]$  such that  $\|\hat{f} - f\|_D \leq \epsilon$ . 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  $\phi : \mathbb{R}^n \times \{0, 1\} \rightarrow [0, 1]$ . In contrast, a query to a generic distribution

$D$  on  $\mathbb{R}^n \rightarrow \mathbb{R}$  takes in real-valued labels, and in Assumption 4.3.1 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 4.2.1.** 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  $\rho_D(\mathcal{C}) = \frac{1}{|\mathcal{C}|} \sum_{c \in \mathcal{C}} \langle c, \mathbb{1}_D \rangle$ . The *statistical dimension on average* at threshold  $\gamma$ ,  $\text{SDA}_D(\mathcal{C}, \gamma)$ , is the largest  $d$  such that for all  $\mathcal{C}' \subseteq \mathcal{C}$  with  $|\mathcal{C}'| \geq |\mathcal{C}|/d$ ,  $\rho_D(\mathcal{C}') \geq \gamma$ .

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

**Theorem 4.2.2** ([GGJ<sup>+</sup>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  $\tau$ , the maximal desired  $L^2$  error is  $\epsilon$ , and that the functions in  $\mathcal{C}$  satisfy  $\|f\|_k \leq \beta$  for all  $f \in \mathcal{C}$ . For technical reasons, we will require  $\tau \leq \epsilon^2$ ,  $\epsilon \leq \beta/3$  (see Section C.1 for some discussion). Then learning  $\mathcal{C}$  up to  $L^2$  error  $\epsilon$  (we may pick  $\epsilon$  as large as  $\beta/3$ ) requires at least  $\text{SDA}_D(\mathcal{C}, \tau^2)$  queries of tolerance  $\tau$ .*

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  $\psi, \phi : \mathbb{R} \rightarrow \mathbb{R}$  and weights  $w_i \in \mathbb{R}^n, a_i \in \mathbb{R}$ .

**Theorem 4.2.3** ([DKKZ20]). *There exists a class  $G$  of one-layer neural networks on  $\mathbb{R}^n$  with  $k$  hidden units such that for some universal constant  $0 < c < 1/2$  and  $\gamma = n^{-k(c-1/2)}$ ,  $\text{SDA}(G, \gamma) \leq 2^{n^c}$ . This holds for any  $\psi : \mathbb{R} \rightarrow [-1, 1]$  that is odd, and  $\phi \in L^2(\mathbb{R}, N(0, 1))$  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 Section C.2.

**Corollary 4.2.4.** *For the following instantiations of  $G$ , with accompanying norm lower bound  $\beta$  (i.e. such that  $\|g\| \geq \beta$  for all  $g \in G$ ), there exist  $\tau = n^{-c}$  and  $\epsilon = \tau$  such that learning  $G$  up to  $L^2$  error  $\epsilon$  requires at least  $2^{n^c}$  queries of tolerance  $\tau$ , for some  $0 < c < 1/2$ .*

- (a) *ReLU nets:  $\psi = \tanh$ ,  $\phi = \text{ReLU}$ . Then  $\beta = \Omega(1/k^6)$  (Lemma C.2.4), so we may take  $\epsilon = \Theta(1/k^6)$ .*
- (b) *Sigmoid nets:  $\psi = \tanh$ ,  $\phi = \sigma$ . Then  $\beta = \exp(-O(\sqrt{k}))$  (Lemma C.2.6), 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  $\epsilon = \Theta(1)$ .*

#### 4.2.2 Convex Optimization Basics

Over a general inner product space  $Z$ , a function  $p : Z \rightarrow \mathbb{R}$  is convex if for all  $\alpha \in [0, 1]$  and  $z, z' \in Z$ ,  $p(\alpha z + (1 - \alpha)z') \leq \alpha p(z) + (1 - \alpha)p(z')$ . We say that  $s \in Z$  is a subgradient of  $p$  at  $z$  if  $p(z + h) \geq p(z) + \langle s, h \rangle$ . We say that  $p$  is  $\beta$ -smoothly convex if for all  $z, h \in Z$  and any subgradient  $s$  of  $p$  at  $z$ ,

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

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  $z, z' \in Z$ ,

$$\langle \nabla p(z), z - z' \rangle \leq \frac{1}{2\beta} \|\nabla p(z) - \nabla p(z')\|^2. \quad (4.2.1)$$

In this chapter 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  $Z^\circ \subset Z$  both be compact, convex subsets of our generic inner product space. Say we have a  $\beta$ -smoothly convex function  $p : Z \rightarrow \mathbb{R}$ , and we want to solve  $\min_{z \in Z^\circ} p(z)$ , i.e. optimize over the smaller domain, while allowing ourselves the freedom of finding subgradients that lie in the larger  $Z$ . 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  $z_0 \in Z$ .

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

    Let  $\gamma_t = \frac{2}{t+2}$ .

    Find  $s \in Z$  such that  $\langle s, -\substack{\text{r} p(z_t) \\ \text{i}} \rangle = \max_{s \in Z^\circ} \langle s, -\substack{\text{r} p(z_t) \\ \text{i}} \rangle - \frac{1}{2} \delta \gamma_t C_p$ .

    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 Section C.4.

**Theorem 4.2.5.** *Let  $Z^\circ \subset Z$  be convex sets, and let  $p : Z \rightarrow \mathbb{R}$  be a  $\beta$ -smoothly convex function. Let  $C_p = \beta \text{diam}(Z)^2$ . For every  $t$ , the iterates of Algorithm 1 satisfy*

$$p(z_t) - \min_{z \in Z^\circ} p(z) \leq \frac{2C_p}{t+2}(1 + \delta).$$

### 4.3 Functional gradient descent

Let  $\ell : \mathbb{R}^n \times \mathbb{R} \rightarrow \mathbb{R}$  be a loss function. Given a  $p$ -concept  $f$  and its corresponding labeled distribution  $D_f$ , the population loss of a function  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  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  $F$  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  $F$  is the convex hull  $\text{conv}(H)$  of a simpler class  $H$ , this can be done using a sufficiently powerful agnostic learning primitive for  $H$ . Thus we can “boost” such a primitive in a black-box manner to minimize  $L(f)$ .

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

**Assumption 4.3.1.** *There is an SQ learner for  $H$  with the following guarantee. Let  $D$  be any labeled distribution on  $\mathbb{R}^n \times \mathbb{R}$  such that the marginal on  $\mathbb{R}^n$  is  $D = N(0, I_n)$ . Given only SQ access to  $D$ , the learner outputs a function  $f \in B(\text{diam}(H)/2)$  such that*

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

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

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

With such an  $H$  in place, we define  $F = \text{conv}(H)$ . We assume that  $f \in F$ . Our objective will be to agnostically learn  $F$ : to solve  $\min_{f \in F} L(f)$  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  $\ell$ .

**Assumption 4.3.2.** *The loss function  $\ell : \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}$  is  $\beta$ -smoothly convex in its first parameter.*

From this assumption, corresponding 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  $f, h : \mathbb{R}^n \rightarrow \mathbb{R}$ . Observe that at for every  $x \in \mathbb{R}^n, y \in \mathbb{R}$ , the subgradient property of  $\ell$  tells us that

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

Taking expectations over  $(x, y) \sim D_f$ , this yields

$$\begin{aligned} L(f + h) - L(f) &= \mathbb{E}_{(x,y) \sim D_f} [\partial_1 \ell(f(x), y)h(x)] \\ &= \mathbb{E}_x [\mathbb{E}_{y|x} [\partial_1 \ell(f(x), y)h(x)]] \\ &= \langle s, h \rangle, \end{aligned}$$

where

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

is thus a subgradient of  $L$  at  $f$ .  $\beta$ -smooth convexity is also easily established. Taking expectations over  $(x, y) \sim D_f$  of the inequality

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

we get

$$L(f + h) - L(f) \geq \langle s, h \rangle - \frac{\beta}{2} \mathbb{E} [h^2]$$

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  $\nabla L(f) : x \mapsto \mathbb{E}_y [\partial_1 \ell(f(x), y)]$ .

**Example 4.3.3.** The canonical example is the squared loss functional, with  $\ell_{\text{sq}}(a, b) = (a - b)^2$ , which is 2-smoothly convex. Here the gradient has a very simple form, since  $\partial_1 \ell_{\text{sq}}(a, b) = 2(a - b)$ , and so

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

i.e.  $\nabla L_{\text{sq}}(f) = 2(f - \mathbb{E}[y])$ . In fact, it is easily calculated that

$$\begin{aligned} L_{\text{sq}}(f) &= \mathbb{E}_{(x,y) \sim D_f} [(f(x) - y)^2] = \mathbb{E}_{(x,y) \sim D_f} [f(x)^2] - 2 \mathbb{E}_{(x,y) \sim D_f} [f(x)y] + \mathbb{E}_{(x,y) \sim D_f} [y^2] \\ &= \mathbb{E}_x [f(x)^2] - 2 \mathbb{E}_x [f(x) \mathbb{E}[y|x]] + \mathbb{E}_{(x,y) \sim D_f} [y^2] \\ &= k f^2 - 2hf, f \text{ i} + 1, \end{aligned}$$

It is also useful to note that

$$L_{\text{sq}}(f) = \int_{\mathcal{X}} L(f(x)) dP(x) = \int_{\mathcal{X}} \frac{1}{2} (f(x) - y)^2 dP(x). \quad (4.3.1)$$

### 4.3.1 Frank–Wolfe using statistical queries

We see that our loss functional is a  $\beta$ -smoothly convex functional on the space  $L^2(\mathcal{R}^n, D)$ . We can now use Frank–Wolfe if we can solve its main subproblem: finding an approximate solution to  $\max_{h \in H} \int_{\mathcal{X}} h(x) \nabla L(f)(x) dP(x)$ , where  $f$  is the current hypothesis during some iteration. Since this is a linear optimization objective and  $F = \text{conv}(H)$ , this is the same as solving  $\max_{h \in H} \int_{\mathcal{X}} h(x) \nabla L(f)(x) dP(x)$ . This is almost the guarantee that Assumption 4.3.1 gives us, but some care is in order. What we have SQ access to is the labeled distribution  $D_f$  on  $\mathcal{R}^n \times \mathcal{Y}$ . It is not clear that we can rewrite the optimization objective in such a way that

$$\max_{h \in H} \int_{\mathcal{X}} h(x) \nabla L(f)(x) dP(x) = \max_{h \in H} \int_{\mathcal{X} \times \mathcal{Y}} h(x) y dD_f(x, y) \quad (4.3.2)$$

for some distribution  $D$  on  $\mathcal{R}^n \times \mathcal{Y}$  that we can simulate SQ access to. Naively, we might try to do this by letting  $D$  be the distribution of  $(x, \nabla L(f)(x))$  for  $x \sim D$ , so that a query  $\phi : \mathcal{R}^n \times \mathcal{Y} \rightarrow \mathcal{R}$  to  $D$  can be answered with  $\int_{\mathcal{X} \times \mathcal{Y}} \phi(x, y) dD(x, y) = \int_{\mathcal{X}} \phi(x, \nabla L(f)(x)) dP(x)$ . 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  $D$ . We turn to the details now.

## 4.4 Functional gradient descent guarantees on surrogate loss

The functional GD approach applied directly to squared loss would allow us to learn  $F = \text{conv}(H)$  using a learner for  $H$  (that satisfied Assumption 4.3.1). But by considering a certain surrogate loss, we can use the same learner to actually learn  $\psi \circ F = \hat{f} \psi$  for an outer activation function  $\psi$ . This is particularly

useful as we can now capture  $p$ -concepts corresponding to functions in  $F$  by using a suitable  $\psi : \mathbb{R} \rightarrow [ -1, 1 ]$ . For example, the common softmax activation corresponds to taking  $\psi = \tanh$ .

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

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

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

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

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

$$\nabla L_{\text{sur}}(f) : x \mapsto \mathbb{E}_{y \sim \psi \circ f(x)} [\partial_1 \ell_{\text{sur}}(f(x), y)] = \psi(f(x)) - \psi(f(x)),$$

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

We still need to show that the Frank–Wolfe subproblem can be solved using

access to just  $D_{\psi \circ f}$ . Observe that

$$\begin{aligned}
\mathbb{E}_x \left[ h(x) \circ L_{\text{sur}}(f)(x) \right] &= \mathbb{E}_x \left[ h(x) (\psi(f(x)) \circ \psi(f(x))) \right] \\
&= \mathbb{E}_x \left[ h(x) \left( \mathbb{E}_{y \sim D_{\psi \circ f}(x)} [y] \circ \psi(f(x)) \right) \right] \\
&= \mathbb{E}_{(x,y) \sim D_{\psi \circ f}} [h(x) (y \circ \psi(f(x)))] \\
&= \mathbb{E}_{(x,y^\theta) \sim D} [h(x) y^\theta],
\end{aligned}$$

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

$$\mathbb{E}_{(x,y^\theta) \sim D} [\phi(x, y^\theta)] = \mathbb{E}_{(x,y) \sim D_{\psi \circ f}} [\phi(x, y \circ \psi(f(x)))] = \mathbb{E}_{(x,y) \sim D_{\psi \circ f}} [\phi^\theta(x, y)] \quad (4.4.1)$$

for the modified query  $\phi^\theta(x, y) = \phi(x, y \circ \psi(f(x)))$ . This means we can rewrite the optimization objective to fit the form in Eq. (4.3.2). Thus for our surrogate loss, Assumption 4.3.1 allows us to solve the Frank–Wolfe subproblem, giving us Algorithm 2 for learning  $F$ .

---

**Algorithm 2** Frank–Wolfe for solving  $\min_{f \in F} L_{\text{sur}}(f)$

---

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

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

    Let  $\gamma_t$  be  $\frac{2}{t+2}$ .

    Let  $D_t$  be the distribution of  $(x, y \circ \psi(f_t(x)))$  for  $(x, y) \sim D_{\psi \circ f}$ .

    Using Assumption 4.3.1, find  $h \in B(\text{diam}(H)/2)$  such that

$$\mathbb{E}_{(x,y^\theta) \sim D_t} [h(x) y^\theta] \geq \max_{h^\theta \in 2H} \mathbb{E}_{(x,y^\theta) \sim D_t} [h^\theta(x) y^\theta] - \frac{1}{2} \gamma_t \lambda \text{diam}(H)^2$$

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

**end for**

---

**Theorem 4.4.1.** *Let  $H$  be a class for which Assumption 4.3.1 holds, and let  $F = \text{conv}(H)$ . Given SQ access to  $D_{\psi \circ f}$  for a known non-decreasing  $\lambda$ -Lipschitz activation  $\psi$  and an unknown  $f \in F$ , suppose we wish to learn  $\psi \circ f$  in terms of surrogate*

loss, i.e. to minimize  $L_{\text{sur}}(f)$ . Then after  $T$  iterations of Algorithm 2, we have the following guarantee:

$$L_{\text{sur}}(f_T) - L_{\text{sur}}(f) \leq \frac{4\lambda \text{diam}(H)^2}{T+2}.$$

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

*Proof.* By the preceding discussion, the surrogate loss functional is  $\lambda$ -smoothly convex, and Algorithm 2 is a valid special case of Algorithm 1, with  $Z = B(\text{diam}(H)/2)$  and  $Z^\theta = \text{conv}(F)$ . Thus the guarantee follows directly from Theorem 4.2.5 (setting  $\delta = 1$ ).

To bound the number of queries, observe that it is sufficient to run for  $T = \frac{4\lambda \text{diam}(H)^2}{\epsilon} + 2$  rounds. In the  $t^{\text{th}}$  iteration, we invoke Assumption 4.3.1 with

$$\epsilon^\theta = \frac{1}{2} \gamma_t \lambda \text{diam}(H)^2 = \frac{\lambda \text{diam}(H)^2}{t+2} - \frac{\lambda \text{diam}(H)^2}{T+2} = \frac{\epsilon}{4}.$$

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

Lastly, we can show that minimizing surrogate loss also minimizes the squared loss. Observe first that  $r L_{\text{sur}}(f) = 0$ . Thus, applying Eq. (4.2.1) with  $z = f$  and  $z^\theta = f$ , we obtain

$$\begin{aligned} L_{\text{sur}}(f) - L_{\text{sur}}(f) &\leq \frac{1}{2\lambda} \|r L_{\text{sur}}(f) - r L_{\text{sur}}(f)\|^2 \\ &= \frac{1}{2\lambda} \|k\psi - f - \psi - f\|^2 \\ &= \frac{1}{2\lambda} (L_{\text{sq}}(\psi - f) + L_{\text{sq}}(\psi - f)), \end{aligned} \tag{4.4.2}$$

where  $L_{\text{sq}}$  is squared loss w.r.t.  $D_{\psi - f}$  and the last equality is Eq. (4.3.1). In particular, Eq. (4.4.2) implies that  $\psi - f$  achieves the following  $L^2$  error with respect to  $\psi - f$ :

$$\|k\psi - f - \psi - f\| \leq \sqrt{2\lambda (L_{\text{sur}}(f) - L_{\text{sur}}(f))}. \tag{4.4.3}$$



## 4.5 Lower bounds on learning ReLUs, sigmoids, and half-spaces

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 4.5.1.** *Let  $H_{\text{ReLU}} = \{x \mapsto \text{ReLU}(w \cdot x) \mid \|w\|_2 \leq 1\}$  be the class of ReLUs on  $\mathbb{R}^n$  with unit weight vectors.<sup>1</sup> Suppose that Assumption 4.3.1 holds for  $H_{\text{ReLU}}$ . Then for any  $\epsilon$ , there exists  $\tau = n^{-c \epsilon^{1/12}}$  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  $G$  from Theorem 4.2.3 instantiated with  $\psi = \tanh$  (which is 1-Lipschitz, so  $\lambda = 1$ ) and  $\phi = \text{ReLU}$ . By the conditions on the weights, we see that  $G \subseteq \tanh \circ F_{\text{ReLU}}$ , where  $F_{\text{ReLU}} = \text{conv}(H_{\text{ReLU}})$ . This construction has a free parameter  $k$ , which we will set based on  $\epsilon$ .
- By our main reduction (Assumption 4.3.1 and Theorem 4.4.1), we can learn  $\tanh \circ F_{\text{ReLU}}$  with respect to  $L_{\text{sur}}$  up to agnostic error  $\epsilon$  using  $O(\frac{1}{\epsilon} q(n, \frac{\epsilon}{4}, \tau))$  queries of tolerance  $\tau$ . By Eq. (4.4.3), this implies learning  $G$  up to  $L^2$  error  $\frac{\rho}{2\epsilon}$ .
- We know that learning  $G$  should be hard. Specifically, Corollary 4.2.4(a) states that if  $\epsilon^\ell = \Theta(1/k^6)$  and the queries are of tolerance  $\tau = n^{-k}$ , then learning up to  $L^2$  error  $\epsilon^\ell$  should require  $2^{n^c}$  queries.

---

<sup>1</sup>We use  $\text{ReLU}$  for simplicity. Any learner can handle this by doing a bit flip on its own.

- The loss our reduction achieves is  $\epsilon^\ell = \frac{\rho_{\overline{2\epsilon}}}{2\epsilon}$ , so we require  $\frac{\rho_{\overline{2\epsilon}}}{2\epsilon} = \Theta(1/k^6)$  for the bound to hold. Accordingly, we pick  $k = \Theta(\epsilon^{-1/12})$ , so that  $\tau = n^{\binom{k}{2}} = n^{\binom{\epsilon^{-1/12}}{2}}$ .
- Thus we must have  $\frac{1}{\epsilon}q(n, \frac{\epsilon}{4}, \tau) \leq 2^{n^c}$ . Rearranging and rescaling  $\epsilon$  gives the result.

□

**Theorem 4.5.2.** *Let  $H_\sigma = \{f(x) = \sigma(w \cdot x) \mid \|w\|_2 \leq 1\}$ , where  $\sigma$  is the standard sigmoid, be the class of sigmoid units on  $\mathbb{R}^n$  with unit weight vectors. Suppose that Assumption 4.3.1 holds for  $H_\sigma$ . Then for any  $\epsilon$ , there exists  $\tau = n^{\binom{\log 1/\epsilon^2}{2}}$  such that  $q(n, \epsilon, \tau) \leq 2^{n^c} \epsilon$  for some  $0 < c < 1/2$ .*

*Proof.* Very similar to the above. We instantiate  $G$  with  $\psi = \tanh$ ,  $\phi = \sigma$ , and observe that  $G = \tanh \circ \text{conv}(H_\sigma)$  and that  $\text{diam}(H_\sigma) = 2$ . In this case, Corollary 4.2.4(b) tells us that we require  $\frac{\rho_{\overline{2\epsilon}}}{2\epsilon} = e^{-\binom{\rho_{\overline{k}}}{k}}$  for the lower bound to hold, so we pick  $k = (\log 1/\epsilon)^2$ . The result now follows exactly as before. □

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 4.3.1 is equivalent. But in Section C.5, 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 4.5.3.** *Let  $H_{\text{hs}} = \{f(x) = \text{sign}(w \cdot x) \mid \|w\|_2 \leq 1\}$  be the class of halfspaces on  $\mathbb{R}^n$  with unit weight vectors. Suppose that Assumption 4.3.1 holds for  $H_{\text{hs}}$ . Then for any  $\epsilon$ , there exists  $\tau = n^{\binom{1/\epsilon}{2}}$  such that  $q(n, \epsilon, \tau) \leq 2^{n^c} \epsilon^3$  for some  $0 < c < 1/2$ .*

*Proof.* To approximate the sign function using a Lipschitz function, we define  $\widetilde{\text{sign}}(x)$  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  $G$  instantiated with  $\psi = \phi = \widetilde{\text{sign}}$  satisfies  $G \subseteq \widetilde{\text{sign}} \cdot \text{conv}(H_{\text{hs}})$ , with  $\text{diam}(G) = 2$ . This is because as noted in Theorem 4.2.3,  $G$  has weights  $a_i \geq 1/k$ , so the sum of halfspaces inside  $\psi$  is always a multiple of  $1/k$ , and  $\widetilde{\text{sign}}$  behaves the same as  $\text{sign}$ .

Theorem 4.4.1 now lets us learn  $G$  up to agnostic error  $\epsilon$  (and hence  $L^2$  error  $\sqrt{2k\epsilon}$ , by Eq. (4.4.3)) using  $O(\frac{k^2}{\epsilon}q(n, \epsilon/4, \tau))$  queries of tolerance  $\tau$ . By Corollary 4.2.4(c), we only need  $\sqrt{2k\epsilon} = \Theta(1)$  for the lower bound to hold, so we may take  $k = \Theta(1/\epsilon)$  to get a lower bound of  $2^{n^c}$ . Thus  $\frac{k^2}{\epsilon}q(n, \epsilon/4, \tau) = 2^{n^c}$ , and rearrangement gives the result.  $\square$

## 4.6 Lower bounds on learning general non-polynomial activations

Here we extend our lower bounds to general non-polynomial activations  $\phi : \mathbb{R} \rightarrow \mathbb{R}$ , by which we mean functions which have an infinite Hermite series  $\phi = \sum_a \hat{\phi}_a H_a$ , where the  $H_a$  are the normalized probabilists' Hermite polynomials. We will again work with the class  $G$  from Theorem 4.2.3, instantiated with this  $\phi$  and  $\psi = \tanh$ . In Section C.2, 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 C.2.1 we show that  $\|kgk\|$  is determined only by  $k$ , the number of hidden units (there  $k = 2m$ ), and the Hermite expansion of  $\phi$ . The reason we require an infinite Hermite series for  $\phi$  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  $\phi$  and  $k$ . In particular, we may assume that it satisfies a norm lower bound  $\|kf\| \geq \beta(k)$ , where  $\beta$  is a function only of  $k$  that is nonzero for infinitely many  $k$ . Here we view the dependence on  $\phi$  as constant.

A few remarks are in order as to how such a bound  $\beta(k)$  may be quantitatively established. If  $\phi$  is either bounded or exhibits only polynomial growth, then the bound on  $\|kgk$  (Lemma C.2.1) gives a corresponding lower bound on  $\|kfk$  that is also purely a function of  $k$ . If  $\phi$  is bounded, the calculation is straightforward and very similar to the  $\phi = \sigma$  case (Lemma C.2.6). If  $\phi$  grows only like a polynomial, then one can use a truncation argument similar to the  $\phi = \text{ReLU}$  case (Lemma C.2.4).

By Theorem 4.2.2 and Corollary 4.2.4, our lower bound of  $2^{n^c}$  on learning  $G$  holds for  $\epsilon \leq \beta(k)/3$ . Since we can pick  $k$  as we like, let us say that for all sufficiently small  $\epsilon$ , we can achieve  $\epsilon \leq \beta(k)/3$  by taking  $k = k(\epsilon) = 3\beta^{-1}(\epsilon)$ . The corresponding tolerance is then  $\tau = n^{-k(\epsilon)}$ , which is still inverse superpolynomial in  $n$ .

We now get the following lower bound on learning  $H = \int \phi(w \cdot x) jkwk_2 \leq 1g$ , again by the same arguments as in Section 4.5. We assume that  $\|k\phi k \leq R$  for some  $R$ , so that  $\text{diam}(H) \leq 2R$ .

**Theorem 4.6.1.** *Suppose that Assumption 4.3.1 holds for  $H$ . Then for all sufficiently small  $\epsilon$  and  $\tau = n^{-k(\epsilon)}$ ,  $q(n, \epsilon, \tau) \geq 2^{n^c} \frac{\epsilon}{R^2}$  for some  $0 < c < 1/2$ .*

*Proof.* We have  $G \subseteq \text{tanh-conv}(H)$ . By functional GD wrt surrogate loss (Theorem 4.4.1), we see that we can learn  $G$  up to  $L^2$  error  $\frac{\rho}{2\epsilon}$  using  $O(\frac{R^2}{\epsilon} q(n, \epsilon, \tau))$  queries of tolerance  $\tau$ , but we must have  $O(\frac{R^2}{\epsilon} q(n, \epsilon, \tau)) \geq 2^{n^c}$ .  $\square$

## 4.7 Lower bounds on learning monomials

In this section we show lower bounds against agnostically learning monomials with respect to the Gaussian, establishing Theorem 4.1.4. Let  $H_{\text{mon}}$  be the class of all multilinear monomials of total degree  $d$  on  $\mathbb{R}^n$ . Clearly  $|H_{\text{mon}}| = \binom{n}{d} = n^{(d)}$ . For any two distinct multilinear monomials  $f, g$ , clearly  $\langle f, g \rangle = 0$  and moreover  $\langle \tanh f, \tanh g \rangle = 0$  as well. Thus the class  $G = \text{tanh-conv}(H_{\text{mon}})$  consists entirely of orthogonal functions. By [GGJ<sup>+</sup>20, Lemma 2.6],  $\text{SDA}(G, \gamma) \geq |G| \gamma = n^{(d)} \gamma$ .

We still need a norm lower bound on  $G$ .

**Lemma 4.7.1.** *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  $\mathbb{E}[\tanh(x_S)^2] \leq \exp(-\Theta(d))$ .*

*Proof.* Observe first that  $\mathbb{E}[x_S^2] = 1$ . By Paley–Zygmund, we have

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

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

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

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

Now since  $\tanh$  is monotonic and odd, we have

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

□

By Theorem 4.2.2 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  $G$  up to  $L^2$  error  $\epsilon$  takes at least  $\text{SDA}(G, \tau^2) \geq n \binom{d}{\tau^2}$  queries.

Now we can use the same arguments as in Section 4.5 to prove the following.

**Theorem 4.7.2.** *Suppose that Assumption 4.3.1 holds for  $H_{\text{mon}}$ . Then for any  $\epsilon \leq \exp(-\Theta(d))$  and  $\tau \leq \epsilon^2$ ,  $q(n, \epsilon, \tau) \geq n \binom{d}{\tau^{5/2}}$ .*

*Proof.* Observe that  $G = \tanh \circ \text{conv}(H_{\text{mon}})$ , and  $\text{diam}(H_{\text{mon}}) \leq 2$ . Using the surrogate loss with  $\psi = \tanh$ , Assumption 4.3.1 and Theorem 4.4.1 tell us that we can learn  $\tanh \circ \text{conv}(H_{\text{mon}})$  up to  $L^2$  error  $\frac{\rho}{2\epsilon}$  (again by Eq. (4.3.1)) in  $O(\frac{1}{\epsilon} q(n, \epsilon, \tau))$  queries of tolerance  $\tau$ . By our lower bound for  $G$ , we must have  $\frac{1}{\epsilon} q(n, \epsilon, \tau) \geq n \binom{d}{\tau^2}$ , or  $q(n, \epsilon, \tau) \geq n \binom{d}{\tau^{5/2}}$  (since  $\epsilon \leq \frac{\rho}{\tau}$ ). □

## 4.8 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  $D = N(0, I_n)$ , the  $\delta$ -approximate degree of a function  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  is the smallest  $d$  such that there exists a degree- $d$  polynomial  $p$  satisfying  $\|f - p\| \leq \delta$ . We show that for any class of  $\delta$ -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  $\epsilon$ , one that satisfies Assumption 4.3.1. We assume bounded labels, say  $y \in [-C, C]$  for some constant  $C$ .

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

Our algorithm will be based on approximating the low-degree Hermite coefficients of  $f_{\text{cmf}}$ , which is equivalent to performing polynomial  $L^2$  regression. It is well-known that in this context, where  $d$  is the  $\delta$ -approximate degree, polynomial  $L^1$  regression up to degree  $d$  gives a squared loss guarantee of  $\delta$  [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  $f_{\text{cmf}}$ ,  $h_{\text{opt}}$ , 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  $h_{\text{opt}}^d$  and  $f_{\text{cmf}}^d$  be the Hermite components of degree at most  $d$  of  $h_{\text{opt}}$  and  $f_{\text{cmf}}$  respectively. Let  $\tilde{f}_{\text{cmf}}^d = \frac{R}{\|f_{\text{cmf}}^d\|} f_{\text{cmf}}^d$ . Among polynomials of degree  $d$  in  $B(R)$ , it is easy to see that  $\tilde{f}_{\text{cmf}}^d$  maximizes  $\langle f, f_{\text{cmf}} \rangle$ , so that

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

Our agnostic learner will look to approximate  $\tilde{f}_{\text{cmf}}^d$  by outputting  $p$  defined as follows. Suppose  $f_{\text{cmf}} = \sum_{I \in 2^{\mathbb{N}^d}} \alpha_I H_I$ , 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  $\beta_I$  be our estimate of  $\alpha_I = \langle f_{\text{cmf}}, H_I \rangle$  to within tolerance  $\tau$  (to be determined). This can be done using  $n^{O(d)}$  queries of tolerance  $\tau$ . Let  $\tilde{f} = \sum_{|I| \leq d} \beta_I H_I$ , and finally let  $p = \frac{R}{k} \tilde{f}$ . We have

$$\begin{aligned}
\left\| \tilde{f}_{\text{cmf}}^d - p \right\|^2 &= R^2 k \left\| \frac{f_{\text{cmf}}^d}{\|f_{\text{cmf}}^d\|} - \frac{\tilde{f}}{\|\tilde{f}\|} \right\|^2 \\
&= R^2 k \left\| \frac{f_{\text{cmf}}^d}{\|f_{\text{cmf}}^d\|} - \frac{\tilde{f}}{\|\tilde{f}\|} \right\|^2 + \left\| \frac{\tilde{f}}{\|\tilde{f}\|} \right\|^2 \left( \frac{1}{\|f_{\text{cmf}}^d\|} - \frac{1}{\|\tilde{f}\|} \right)^2 \\
&= 2R^2 \left( \frac{\|f_{\text{cmf}}^d - \tilde{f}\|^2}{\|f_{\text{cmf}}^d\|^2} + \|\tilde{f}\|^2 \left( \frac{1}{\|f_{\text{cmf}}^d\|} - \frac{1}{\|\tilde{f}\|} \right)^2 \right) \\
&= 2R^2 \left( \frac{\|f_{\text{cmf}}^d - \tilde{f}\|^2}{\|f_{\text{cmf}}^d\|^2} + \left( \frac{\|f_{\text{cmf}}^d\| - \|\tilde{f}\|}{\|f_{\text{cmf}}^d\|} \right)^2 \right) \\
&= 4R^2 \frac{\|f_{\text{cmf}}^d - \tilde{f}\|^2}{\|f_{\text{cmf}}^d\|^2} \quad (\text{triangle ineq.}) \\
&= \frac{4R^2 n^d \tau^2}{k \|f_{\text{cmf}}^d\|^2}, \tag{4.8.1}
\end{aligned}$$

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

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

$$Rk \langle f_{\text{cmf}}^d, h_{\text{opt}} \rangle = \langle f_{\text{cmf}}^d, h_{\text{opt}} \rangle \leq \|f_{\text{cmf}}^d\| \|h_{\text{opt}}\| \leq \|f_{\text{cmf}}^d\| k \delta,$$

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

By Eq. (4.8.1), we then have

$$\| \tilde{f}_{\text{cmf}}^d - p \| \leq \frac{4Rn^{d/2}\tau}{\epsilon}. \quad (4.8.2)$$

Now observe that

$$\begin{aligned} | \langle hp, f_{\text{cmf}} \rangle - \langle h\tilde{f}_{\text{cmf}}^d, f_{\text{cmf}} \rangle | &= | \langle h\tilde{f}_{\text{cmf}}^d, f_{\text{cmf}} \rangle - \langle hp, f_{\text{cmf}} \rangle | \\ &\leq \frac{4RCn^{d/2}\tau}{\epsilon} \quad (\text{Eq. (4.8.2) and Cauchy-Schwarz}) \\ &= | \langle hh_{\text{opt}}^d, f_{\text{cmf}} \rangle - \langle hp, f_{\text{cmf}} \rangle | + | \langle hp, f_{\text{cmf}} \rangle - \langle hh_{\text{opt}}^d, f_{\text{cmf}} \rangle | \\ &\leq \frac{\epsilon}{2} + \frac{4RCn^{d/2}\tau}{\epsilon}. \quad (\text{Cauchy-Schwarz, and using } \delta C = \epsilon/2) \end{aligned}$$

Setting  $\tau = \frac{\epsilon^2}{8RCn^{d/2}}$  gives us the desired result, namely that  $| \langle hp, f_{\text{cmf}} \rangle - \langle hh_{\text{opt}}^d, f_{\text{cmf}} \rangle | \leq \epsilon$ .

Thus we have the following theorem.

**Theorem 4.8.1.** *The class  $H_{\text{ReLU}}$  can be agnostically learned up to correlation  $\epsilon$  (in the sense of Assumption 4.3.1) using  $n^{O(\epsilon^{-4/3})}$  queries of tolerance  $n^{-\epsilon^{-4/3}}\epsilon$ . Similarly,  $H_\sigma$  can be learned using  $n^{O(\log^2 1/\epsilon)}$  queries of tolerance  $n^{-(\log^2 1/\epsilon)}\epsilon^2$ .*

*Proof.* Approximating the Hermite coefficients of degree at most  $d$  takes  $n^{O(d)}$  queries of tolerance  $n^{-d}\epsilon$ . As we show in Section C.3, the  $\delta$ -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.  $\square$

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



# Chapter 5: The Polynomial Method is Universal for Distribution-Free Correlational SQ Learning

## 5.1 Introduction

A successful and general approach to distribution-free learning in both Valiant’s PAC model [Val84] and Kearns *et al.*’s agnostic model [KSS94] has been the so-called “polynomial method.” In this approach, a Boolean function is approximated by a low-degree polynomial with respect to various losses such as the 0/1 loss of the sign of the polynomial or a uniform approximation by the polynomial itself. This representation in turn leads to algorithms whose complexity is typically exponential in the degree of the approximating polynomial.

Several researchers have pointed out that the polynomial method captures the best-known distribution-free learners for Boolean function classes with the exception of classes that include parities [HS07]. It is used for example in the algorithm of Klivans and Servedio [KS04] for PAC-learning DNF formulas as well as the  $L^1$  polynomial regression algorithm of Kalai *et al.* [KKMS08] for agnostically learning conjunctions.

In this work we give a complete characterization of (correlational) statistical query learning in the distribution-free PAC or agnostic setting in terms of the threshold degree or approximate degree of the unknown function class. This shows that unless we can find new algorithms that lie outside the SQ model, our only approach for distribution-free learning is to construct low-degree approximating polynomials.

To be more concrete, let  $f : \{-1, 1\}^n \rightarrow \{-1, 1\}$  be a Boolean function. The *threshold degree*  $\deg(f)$  of  $f$  is the least degree of a polynomial  $p$  such that  $f(x) = \text{sign}(p(x))$  for all  $x$ . The *pointwise  $\epsilon$ -approximate degree*  $\widetilde{\deg}_\epsilon(f)$  of  $f$  is the least degree of a polynomial  $p$  such that  $|f(x) - p(x)| \leq \epsilon$  for all  $x$ . It is well-known that for any class  $\mathcal{C}$  of functions with threshold degree bounded by  $d$ ,  $\mathcal{C}$  can

be viewed as a class of halfspaces over an  $n^{O(d)}$ -dimensional space (corresponding to the monomials up to degree  $d$ ), and classic halfspace learners then yield distribution-free PAC learners for  $\mathcal{C}$  that run in  $n^{O(d)}$  time. Similarly in the agnostic setting, the results of [KKMS08] imply (as noted for example in [KS10, Prop 2.1]) that for a class  $\mathcal{C}$  with  $\epsilon$ -approximate degree bounded by  $d$ , degree- $d$   $L^1$  polynomial regression yields a distribution-free agnostic learner (with error  $\text{opt} + \epsilon$ ) running in time  $n^{O(d)}$ . This is since pointwise or  $L^1$  approximation implies  $L^2$  approximation (as required by [KKMS08]) with respect to all distributions.

Our main results can now be stated as follows. Formal definitions may be found in the preliminaries.

**Theorem 5.1.1.** *Let  $\mathcal{C}$  be a Boolean function class satisfying some mild conditions (namely being closed under pattern restrictions, see Definition 5.2.12), with threshold degree  $\Omega(d)$ . Any distribution-free PAC learner for  $\mathcal{C}$  using only correlational statistical queries of tolerance  $\tau < \frac{1}{10}$  requires at least  $2^{\binom{d}{2}\tau^2}$  queries in order to learn  $\mathcal{C}$  up to error  $\frac{1}{3}$ .*

**Theorem 5.1.2.** *Let  $\mathcal{C}$  be a Boolean function class as above (closed under pattern restrictions), with  $\frac{1}{2}$ -approximate degree  $\Omega(d)$ . Any distribution-free agnostic learner for  $\mathcal{C}$  using only correlational statistical queries of tolerance  $\tau < \frac{1}{10}$  requires at least  $2^{\binom{d}{2}\tau^2}$  queries in order to agnostically learn  $\mathcal{C}$  up to excess error  $\frac{1}{100}$ , i.e. true error  $\text{opt} + \frac{1}{100}$ .*

Both these lower bounds match (up to logarithmic factors in the exponent) the upper bounds given by the polynomial method that were mentioned earlier.

As one example of a corollary, we obtain essentially tight lower bounds for the important problem of agnostically learning conjunctions (likewise disjunctions). Specifically, it is a classic result [NS94, Pat92] that the class of conjunctions on  $f \in \{0, 1\}^n$  has  $\epsilon$ -approximate degree  $\Theta\left(\binom{n}{\epsilon}\right)$  for all constant  $\epsilon < 1$ , and this yields a  $2^{\binom{n}{\epsilon}}$  correlational SQ lower bound by the above theorem.

We point out that a slightly different characterization of CSQ PAC learnability in terms of the existence of low-weight threshold representations had been given in [Fel08] (see Theorem 5.4 therein). Combining this with work by Sherstov [She08b, She11] on the relationships between various complexity measures of Boolean function classes as well as the links between learning and communication complexity, one could already obtain a qualitatively equivalent version of Theorem 5.1.1. Similarly but less directly, the ideas used to obtain Theorem 5.1.2 are arguably implicit in [Fel12] in combination with [She08b, She11]. Thus the primary appeal of our results lies less in their novelty and more in the following factors: (a) the statements are clean and explicitly reveal the dependence on both the degree  $d$  and the tolerance  $\tau$ ; (b) the proofs are simple, largely self-contained (without reference to communication complexity results), and explicitly state the hard instances for which the lower bound is obtained.

Both our theorems are based on a more general theorem that constructs a hard family of functions from a single function  $f$  and a distribution  $\mu$  which “orthogonalizes” it up to degree  $d$ , i.e. under which all degree- $d$  or lower moments of  $f$  vanish. This condition has been used in prior works for proving lower bounds for learning neural networks with respect to Gaussian distributions (e.g. [DKKZ20, DKZ20a]).

This general theorem crystallizes a basic principle of SQ lower bounds, saying that in the SQ setting, the hard distributions for a class of functions are the ones that orthogonalize them, i.e. zero out their low-degree moments:

**Theorem 5.1.3.** *Let  $f : \{-1, 1\}^{n/2} \rightarrow \mathbb{R}$  be a Boolean function, and let  $\mu$  be a distribution that orthogonalizes  $f$  up to degree  $d - 1$ . There exists a family  $A$  of function-distribution pairs on  $\{-1, 1\}^n$  where each function in the family is of the form  $x \mapsto f(x_V \oplus w)$  for some  $V \subseteq [n]$  of size  $n/2$  and  $w \in \{-1, 1\}^{n/2}$ , such that any distribution-free PAC learner for  $A$  using only correlational statistical queries of tolerance  $\tau < \frac{1}{10}$  requires at least  $\Omega(2^d \tau^2)$  queries in order to learn  $C$  up to error  $\frac{1}{3}$ . (Here  $x_V$  denotes the vector of  $x_i$  for  $i \in V$ , and  $\oplus$  denotes the bitwise XOR.)*

The proof is a simple application of the pattern matrix method due to Sherstov [She11]. As a partial converse, it is not hard to see that distributions under which the target functions have significant low-degree moments allow Fourier-theoretic methods (e.g. the “low degree algorithm”) to succeed. We view this result as providing an additional, more general sense in which the polynomial method is a “universal” approach to learning.

One unusual feature of our results is that rather than constructing a single hard distribution under which a large class of functions is hard to learn, our hard families involve a set of function-distribution pairs, where each function is effectively accompanied by its own hard distribution. A distribution-free learner should, of course, succeed equally well on such a family.

Our results also have consequences for the hardness of approximation by linear classes of classes with high threshold degree.

**Theorem 5.1.4.** *Let  $C$  be a Boolean function class closed under pattern restrictions, with threshold degree  $\Omega(d)$ . Consider the linear hypothesis class with an arbitrary  $N$ -dimensional embedding  $\psi$  and a norm bound  $B$ :*

$$H_{\psi,B} = \{x \mapsto \psi(x) \cdot w \mid \|w\|_2 \leq B\}.$$

*There exists a family  $A$  of function-distribution pairs, with the functions lying in  $C$ , such that*

$$\max_{(f,D) \in A} \min_{h \in H_{\psi,B}} \mathbb{E}_{x \sim D} [(f(x) - h(x))^2] > \frac{1}{2} \left( \frac{B}{N} \right)^{\Omega(d)}.$$

As one application, we see that for all classes with threshold degree  $d = \omega(\log n)$  (which includes most classes, e.g. conjunctions, intersections of halfspaces, and  $\text{AC}^0$ ), so that  $2^d$  is superpolynomial, even weak approximation is impossible unless  $B \leq N^{2^{-\Omega(d)}}$ . In particular, weak approximation is impossible with  $B, N = \text{poly}(n)$ . This generalizes earlier results proved for classes such as conjunctions via the approximate rank [KS10].

Our approach builds closely on the recent work of [MS22], in which the authors introduce a new measure called the “variance” of a class of labeled distributions, and use it to show CSQ lower bounds as well as a host of hardness of approximation results for the class. In this chapter we use the term “correlational variance” to refer to this quantity for clarity. They use this new framework to show an essentially-tight  $2^{\Omega(n^{1/3})}$  CSQ lower bound for learning DNF formulas by applying technical results of Razborov and Sherstov [RS10]. Our main contribution is to simplify their approach and extend this framework to the more general setting of orthogonalizing distributions. We can then leverage known lower bounds on the threshold and approximate degree of function classes. For example, we can obtain the  $2^{\Omega(n^{1/3})}$  lower bound of [MS22] by using our main theorem with the classical  $\Omega(n^{1/3})$  lower bound due to Minsky and Papert from the 60s [MP69]. This deviates from established frameworks for proving CSQ lower bounds, which generally use variants of the SQ dimension (see [Fel16, Rey20] for surveys).

### 5.1.1 Related and prior work

There has been a wealth of work on lower bounds for learning various Boolean function classes in various settings, and we can only hope to survey a small slice of it. All Boolean functions in what follow will be assumed to be on  $n$  bits and  $\text{poly}(n)$ -sized.

**Cryptographic lower bounds** One of the original lines of work on the hardness of learning was based on cryptographic assumptions (see [KV94b, Chap. 6] for a textbook reference). Valiant in his original paper [Val84] observed that pseudorandom families, and the function classes that can compute them, are inherently hard to learn in polynomial time. This was developed further by [KV94a], who proved hardness of polynomial-time learning for Boolean formulas and finite automata based on assumptions underlying public-key cryptography. Another important early work was [Kha93], which showed a quasipolynomial time lower bound for learning  $\text{AC}^0$

even under the uniform distribution (matching [LMN93]) based on secure pseudo-random generators. For intersections of halfspaces, [KS09] showed the hardness of polynomial-time learning assuming the security of lattice-based cryptosystems.

**Lower bounds on halfspace-based methods** Historically one of the chief approaches to learning various classes has been the use of halfspace-based learners such as Perceptron, SVM, and linear programming, including kernel methods and in particular the polynomial method. Measures such as approximate rank [BdW01], dimension and margin complexity [Vap00, CST<sup>+</sup>00, BDES02, LMSS07] and their very recent probabilistic variants [KMS20] are known to characterize inherent limitations of this approach. In [KS10], an approximate rank bound of  $2^{\binom{p}{n}}$  was shown for disjunctions, and in [RS10] a dimension complexity bound of  $2^{(n^{1/3})}$  was shown for DNFs. Recent nearly-optimal lower bounds on threshold and approximate degree [She13, BT19, SW19] for classes such as intersections of halfspaces and  $AC^0$ , as well as degree-weight tradeoffs [STT12], may also be considered to fall into this line of work.

**Complexity-theoretic lower bounds** In more recent times there has been a line of work [DLSS14, Dan16b, DS16] establishing hardness of distribution-free learning using complexity theoretic assumptions such as the hardness of refuting random  $k$ -SAT. These works establish hardness of polynomial-time learning for DNF formulas and intersections of  $\omega(\log n)$  halfspaces, and hardness of agnostic learning for conjunctions, halfspaces, and parities.

**SQ lower bounds** The SQ model is the most general one in which unconditional lower bounds have been shown for many classes. The canonical result here is the SQ lower bound of  $2^{\binom{n}{k}}$  for learning parities under the uniform distribution [Kea98, BFJ<sup>+</sup>94], which introduced the SQ dimension. This measure has since been considerably generalized [Fel12, FGR<sup>+</sup>17, Fel17]. As mentioned in the introduction,

other important works in this overall line are those of [Fel08, Fel12], which provided characterizations and lower bounds for distribution-free CSQ learnability in both the PAC and agnostic settings. These can in turn be related to other commonly studied complexity measures using the far-reaching results of [She08b, She11]. Our approach builds on the slightly different work of [MS22], which proved SQ lower bounds in terms of a measure called variance, variants of which had been considered before in [SSSS17, Sha18b].

In the distribution-specific agnostic setting, the polynomial method takes the form of  $L^1$  polynomial regression, introduced by [KKMS08] as a universal approach to agnostic learning. The results of [DSFT<sup>+</sup>14, DKPZ21] later showed matching SQ lower bounds in terms of the  $L^1$  approximate degree of a concept class, establishing this quantity as an appropriate measure of agnostic learnability (at least over the uniform distribution on the hypercube or the standard Gaussian). Our results may be seen as analogs of these results for distribution-free PAC and agnostic learnability.

**SQ vs CSQ** The relationship between the general SQ model and the CSQ restriction has been a topic of much interest. In the distribution-specific setting, the two are unconditionally equivalent [BF02]. More generally, the two are known to be equivalent when either the distribution is known, or sampling or nonuniformity is allowed [Fel08]. But somewhat surprisingly, the two are not in general equivalent in the distribution-free setting: [Fel08] showed that halfspaces are not efficiently CSQ-learnable even though they are always SQ-learnable. Thus our CSQ lower bounds do not readily extend to the general SQ model, and a very compelling direction for future work is to investigate whether such an extension is possible.

**The polynomial method and known upper bounds** Essentially the only known distribution-free approach to learning most rich Boolean classes remains the polynomial method, i.e. representing Boolean functions as polynomial threshold functions (see [HS07] for a survey) or, for agnostic learning, pointwise approximation (which in

particular implies  $L^1$  approximation [KKMS08] with respect to any distribution). A notable exception in the perfectly realizable case is the use of (linear) algebraic methods for learning classes such as parities or more generally polynomials over finite fields [HS07]. It is significant that the former category of methods may all be implemented using simple linear programming and made to work in the SQ setting, whereas the latter category is the prime (and essentially only) example of PAC algorithms that are not SQ.

Improving upon the polynomial method has long been an elusive goal for many important classes. Indeed, in [ST17] the authors propose a change of perspective where rather than seek an efficient running time, the goal of a distribution-free learner is merely to run with nontrivial savings over exponential time, namely in time  $2^{n - s(n)}$  where  $s(n) = \omega(\log n)$ . For  $\text{AC}^0$  circuits of size  $M$  and depth  $d$ , they give such an algorithm with running time  $2^{n - (n/(\log M)^{d-1})}$ .

## 5.2 Preliminaries

**Notation** Let  $f : \{0, 1\}^n \rightarrow \{0, 1\}$  be a Boolean function and  $D$  be a distribution on  $\{0, 1\}^n$ . We will sometimes view such functions as members of the  $L^2$  space  $L^2(\{0, 1\}^n, D)$ , with the inner product given by  $\langle f, g \rangle_D = \mathbb{E}_D[fg]$ . When we use just  $\langle f, g \rangle$  without any subscript, we will mean the inner product with respect to the uniform distribution on  $\{0, 1\}^n$ . We denote by  $f(D)$  the labeled distribution (on  $\{0, 1\}^n \times \{0, 1\}$ ) of  $(x, f(x))$  for  $x \sim D$ . We will generally denote unlabeled distributions on  $\{0, 1\}^n$  by non-calligraphic letters such as  $D$  or the Greek  $\mu$ , and labeled distributions on  $\{0, 1\}^n \times \{0, 1\}$  by calligraphic  $D$  (unless it can be described as  $f(D)$  for some  $f, D$ ).

We use  $[n]$  to denote  $\{1, \dots, n\}$ . Given  $x \in \{0, 1\}^n$  and  $V \subseteq [n]$ , we denote by  $x_V$  the vector of  $x_i$  for  $i \in V$ . Given  $w \in \{0, 1\}^n$ , we denote by  $x \oplus w$  the bitwise XOR (or in our case elementwise product, since we represent bits by  $\{0, 1\}$ ) of  $x$  and  $w$ .



We also make use of basic notions from Boolean Fourier analysis. We use  $\chi_S$  to denote the parity on  $S \subseteq [n]$ ,  $\chi_S(x) = \prod_{i \in S} x_i$ . The Fourier coefficients of a function  $f : \{0, 1\}^n \rightarrow \mathbb{R}$  are denoted  $\hat{f}(S)$ , with  $\hat{f}(S) = \langle f, \chi_S \rangle$ , and  $f = \sum_{S \subseteq [n]} \hat{f}(S) \chi_S$ .

### 5.2.1 Learning in the statistical query model

The statistical query (SQ) model was introduced by Kearns [Kea98] and has proved highly influential as a realistic learning model that also allows strong lower bounds; see [Fel16, Rey20] for surveys. Let  $D$  be a labeled distribution on  $\{0, 1\}^n \times \{0, 1\}$ . An SQ oracle for  $D$  is one that takes as input a query function  $\phi : \{0, 1\}^n \rightarrow [-1, 1]$  and a tolerance  $\tau \in (0, 1)$ , and responds with a value in  $[\mathbb{E}_D[\phi] - \tau, \mathbb{E}_D[\phi] + \tau]$ . When a query is of the special form  $\phi(x, y) = g(x)y$  for some  $g : \{0, 1\}^n \rightarrow [-1, 1]$ , it is known as a correlational query, and is fully specified just by the function  $g$ . In the common case where  $D$  is actually  $f(D)$  for some function  $f : \{0, 1\}^n \rightarrow \{0, 1\}$  and distribution  $D$  on  $\{0, 1\}^n$ , the expectation of a correlational query is  $\mathbb{E}_{(x,y) \sim f(D)}[g(x)y] = \mathbb{E}_x \mathbb{E}_y [f(x)g(x)y] = \langle f, g \rangle_D$ , and for this reason such queries are also called inner product queries.

We recall the definitions of (realizable) PAC and agnostic learning as applicable in the SQ model. Let  $\mathcal{C}$  be a Boolean function class on  $\{0, 1\}^n$ . In the traditional, realizable PAC setting, the learner is given SQ access to  $c(D)$  for an unknown  $c \in \mathcal{C}$  and arbitrary distribution  $D$  on  $\{0, 1\}^n$ , and is said to learn  $\mathcal{C}$  up to error  $\epsilon$  if it is able to output a function  $h$  such that  $\mathbb{P}_{x \sim D}[h(x) \neq c(x)] \leq \epsilon$ . The learner is called a distribution-free learner if it has this guarantee irrespective of what  $D$  is. In the agnostic setting, the labels no longer arise from a function  $c \in \mathcal{C}$ . Instead the learner is given SQ access to an arbitrary labeled distribution  $D$  on  $\{0, 1\}^n \times \{0, 1\}$ , and the goal is to be competitive with the best-fitting classifier in  $\mathcal{C}$ . Letting  $\text{opt} = \inf_{c \in \mathcal{C}} \mathbb{P}_{(x,y) \sim D}[c(x) \neq y]$ , a learner is said to agnostically learn  $\mathcal{C}$  up to (excess) error  $\epsilon$  if it is able to output a function  $h$  such that  $\mathbb{P}_{(x,y) \sim D}[h(x) \neq y] \leq \text{opt} + \epsilon$ . Again, the learner is called distribution-free if this holds for all  $D$ .

For distribution-specific learners of Boolean classes, it is a well-known observation [BF02] that correlational SQs (CSQs) are equivalent to general SQs. This is not known to be the case in the distribution-free setting in which we operate, but it is known that a general SQ learner implies a CSQ learner when either sampling or nonuniformity is allowed [Fel08].

### 5.2.2 Correlational variance

In [MS22], the authors introduce a measure called the “variance” of a function class, or more precisely a family of labeled distributions, and use it to show lower bounds both on CSQ learning as well as approximation of this class. We use the term “correlational variance” to refer to this quantity for clarity. The setting is as follows. Let  $A$  be a family of function-distribution pairs  $(f, D)$ , where  $f : \{0, 1\}^n \rightarrow \{0, 1\}$  and  $D$  is a distribution on  $\{0, 1\}^n$ .  $A$  can also be seen as the family of the labeled distributions  $f(D)$ . Such a family will play exactly the same role as the notion of the “assumption class” from [KSS94]: namely, it is from an unknown member of this family that the learner receives labeled data (or SQ access to it).

**Definition 5.2.1** (Correlational Variance). Let  $A$  be a family of function-distribution pairs as above. For any  $\phi : \{0, 1\}^n \rightarrow [-1, 1]$ , define

$$\text{Var}(A, \phi) = \mathbb{E}_{(f,D) \in A} \langle f, \phi \rangle_D^2 = \frac{1}{|A|} \sum_{(f,D) \in A} \langle f, \phi \rangle_D^2.$$

Then the “correlational variance” of  $A$  is defined as

$$\text{Var}(A) = \sup_{\phi \in \{-1, 1\}^n} \text{Var}(A, \phi).$$

In words, the correlational variance of a family captures the maximum variance (really the second moment) of the correlation of a query function  $\phi$  with a random member of the family.

This measure yields lower bounds on SQ learning using correlational queries. The core of the proof is the following appealingly simple Markov/Chebyshev-style

lemma, which bounds the number of functions in a family that can be highly correlated with a query. We include it with a short proof so as to capture the essential function of the definition of correlational variance.

**Lemma 5.2.2** ([MS22], Lemma 6). *Let  $A$  be a family as above. Fix a query function  $\phi : \mathcal{F} \rightarrow [0, 1]^n$  and a tolerance  $\tau$ . Let  $A_\phi$  denote the function-distribution pairs with correlation at least  $\tau$  with  $\phi$ , i.e.  $\langle \phi, f \rangle_D \geq \tau$ . Then  $|A_\phi| \leq \frac{\text{Var}(A)}{\tau^2} |A|$ .*

*Proof.* What fraction of pairs  $(f, D)$  in  $A$  satisfy  $\langle \phi, f \rangle_D \geq \tau$ ? By a simple Markov/Chebyshev bound, this fraction is at most  $\frac{\text{Var}(A, \phi)}{\tau^2} = \frac{\text{Var}(A)}{\tau^2}$ .  $\square$

Roughly speaking, this means that an adversarial SQ oracle that responds with 0 to every query only allows the learner to rule out at most  $\frac{\text{Var}(A)}{\tau^2} |A|$  functions per query, essentially forcing even a weak learner to use  $\Omega(\frac{\tau^2}{\text{Var}(A)})$  queries in total. Formally, [MS22] show the following bound.

**Theorem 5.2.3** ([MS22], Theorem 7). *Let  $A$  be a family of function-distribution pairs as above. Consider a learner given CSQ access to an unknown  $(f, D) \in A$  with the goal of finding  $h$  such that the 0-1 error  $\mathbb{P}_{x \sim D}[f(x) \neq h(x)]$  is small. Then any such learner making only correlational queries of tolerance  $\tau$  must use  $\Omega(\frac{\tau^2}{\text{Var}(A)})$  queries in order to output a function with 0-1 error better than  $\frac{1}{2} - \frac{\tau}{2}$ . In particular, assuming  $\tau = \frac{1}{10}$  (say),  $\Omega(\frac{\tau^2}{\text{Var}(A)})$  queries are required in order to obtain 0-1 error at most  $\frac{1}{3}$ .*

The correlational variance of a family  $A$  can be bounded in terms of the spectral norm of the linear operator  $M(A)$  mapping  $\phi : \mathcal{X} \rightarrow [0, 1]$  to the vector  $(\langle \phi, f \rangle_D)_{(f, D) \in A}$ . Here the domain  $\mathcal{X}$  will be  $\mathcal{F}$  for us. If we vectorize  $\phi$  as  $v(\phi) = (\phi(x))_{x \in \mathcal{X}}$ , then  $M(A)$  is the matrix of size  $|A| \times |\mathcal{X}|$  whose rows are indexed by  $(f, D) \in A$  and the columns by  $x \in \mathcal{X}$ , and whose entries are given by  $[f(x)D(x)]_{(f, D) \in A, x \in \mathcal{X}}$ . We then have the following bound.

**Lemma 5.2.4.**

$$\text{Var}(A) = \frac{jXj}{jAj} kM(A)k^2.$$

*Proof.* This follows immediately from the observation that for any  $\phi : X \rightarrow [ -1, 1]$ ,

$$\text{Var}(A, \phi) = \frac{1}{jAj} kM(A)v(\phi)k^2 = \frac{1}{jAj} kM(A)k^2 kv(\phi)k^2 = \frac{jXj}{jAj} kM(A)k^2,$$

where  $kv(\phi)k^2 = jXj$  since  $\phi : X \rightarrow [ -1, 1]$ . □

Bounds on correlational variance also have implications for hardness of approximation, which we cover in Section 5.6.

### 5.2.3 Orthogonalizing distributions

The notion of an orthogonalizing distribution will be important to us, as it is the most general setting in which our results can be stated. It is a notion that has been used in many prior CSQ bounds and communication complexity results [She08a, She11, RS10, DKKZ20, DKZ20a].

**Definition 5.2.5** (Orthogonalizing distribution). Let  $f : \{ -1, 1\}^n \rightarrow \{ -1, 1\}$  be a Boolean function. We say a distribution  $\mu$  on  $\{ -1, 1\}^n$  orthogonalizes  $f$  up to degree  $d$  if for all polynomials of degree at most  $d$ ,  $hf, p_i_\mu = 0$ . In particular,  $hf, \chi_S i_\mu = 0$  for all  $|S| \leq d$ .

For us a convenient way to use this definition will be to define a function  $g = f\mu$ , and to observe that for all  $|S| \leq d$ ,

$$\widehat{g}(S) = hg, \chi_S i = 2^{-n} \sum_{x \in \{ -1, 1\}^n} f(x)\mu(x)\chi_S(x) = 2^{-n} hf, \chi_S i_\mu = 0,$$

as well as that for all  $|S| > d$ , since  $hf, \chi_S i_\mu = 0$ ,

$$\widehat{g}(S) = hg, \chi_S i = 2^{-n} hf, \chi_S i_\mu = 0.$$

### 5.2.4 Threshold and approximate degree

Two of the most basic ways of using real polynomials to represent a Boolean function are to represent it as a polynomial threshold function, and to approximate it pointwise. These are both classical notions that have a long history in approximation theory, and both notions are accompanied by a beautiful duality theory. In both cases the dual characterization can be viewed in terms of orthogonalizing distributions.

**Definition 5.2.6** (Threshold degree). Let  $f : \{-1, 1\}^n \rightarrow \{-1, 1\}$  be a function. The threshold degree of  $f$ , denoted  $\deg_\theta(f)$ , is the least degree of a real polynomial  $p : X \rightarrow \mathbb{R}$  such that  $f(x) = \text{sign}(p(x))$  for all  $x$ .

It is a classical result that the threshold degree has a dual characterization in terms of an orthogonalizing distribution; see e.g. [She11, Theorem 3.3].

**Theorem 5.2.7.** *Let  $f : \{-1, 1\}^n \rightarrow \{-1, 1\}$  be a function with threshold degree  $d$ . Then there exists a distribution  $\mu$  on  $X$  that orthogonalizes  $f$  up to degree  $d - 1$ .*

**Definition 5.2.8** (Approximate degree). Let  $f : \{-1, 1\}^n \rightarrow \{-1, 1\}$  be a function. The pointwise  $\epsilon$ -approximate degree of  $f$ , denoted  $\widetilde{\deg}_\epsilon(f)$ , is the least degree of a real polynomial  $p : \{-1, 1\}^n \rightarrow \mathbb{R}$  such that  $|f(x) - p(x)| \leq \epsilon$  for all  $x$ .

Again, the following dual characterization is well-known; see e.g. [She11, Theorem 3.2].

**Theorem 5.2.9.** *Let  $f : \{-1, 1\}^n \rightarrow \{-1, 1\}$  be a function with  $\widetilde{\deg}_\epsilon(f) = d$ . Then there exists a function  $\psi : \{-1, 1\}^n \rightarrow \mathbb{R}$  such that:*

- (a)  $\sum_x \psi(x) f(x) > \epsilon$ ,
- (b)  $\sum_x \psi(x) = 1$ ,
- (c)  $\int \psi p = 0$  for all polynomials  $p$  of degree less than  $d$ .

In particular, if we define a distribution  $\mu(x) = \psi(x)$  and let  $h(x) = \text{sign}(\psi(x))$ , then we see that  $\mu$  orthogonalizes  $h$  up to degree  $d - 1$ , and that  $\int h f d\mu > \epsilon$ .

### 5.2.5 Pattern matrices

Our constructions of hard families make use of Sherstov’s pattern matrix method [She11], along the lines of the construction of [MS22] for DNFs.

**Definition 5.2.10** (Pattern matrix). Let  $f : \{0, 1\}^k \rightarrow \mathbb{R}$  be a function, and let  $n > k$  be a multiple of  $k$ . Let  $V(n, k)$  be the family of all size- $k$  subsets  $V \subseteq [n]$  of the following form: dividing  $[n]$  into  $k$  consecutive blocks of size  $n/k$ ,  $V$  consists of exactly one index from each block. (Notice that  $|V(n, k)| = \binom{n/k}{k}$ .) The  $(n, k, f)$ -pattern matrix is the matrix of size  $2^n \times \binom{n/k}{k} 2^k$  whose rows are indexed by  $x \in \{0, 1\}^n$  and columns by pairs  $(V, w) \in V(n, k) \times \{0, 1\}^k$ , and whose elements are given by  $f(x_V \oplus w)$ .

The following definitions are similar but will make many of our theorems easier to state.

**Definition 5.2.11** (Pattern matrix class). Let  $f : \{0, 1\}^k \rightarrow \mathbb{R}$  be a function, and let  $n > k$  be a multiple of  $k$ . The  $(n, k, f)$ -pattern matrix class is the class of all functions on  $\{0, 1\}^n$  of the form  $x \mapsto f(x_V \oplus w)$  for some  $(V, w) \in V(n, k) \times \{0, 1\}^k$ . We also sometimes use the term  $(n, k, f)$ -pattern matrix family for a family of function-distribution pairs where the functions involved are generated from  $f$  in this way.

**Definition 5.2.12** (Pattern restrictions). Let  $\mathcal{C} = \bigcup_{n>0} \mathcal{C}_n$  be the union of some classes  $\mathcal{C}_n$  of Boolean functions on  $\{0, 1\}^n$ . We say  $\mathcal{C}$  is closed under pattern restrictions if for any  $k, n$  a multiple of  $k$ , and any  $f \in \mathcal{C}_k$ , the function  $x \mapsto f(x_V \oplus w)$  on  $\{0, 1\}^n$  lies in  $\mathcal{C}_n$  for any  $V \subseteq [n]$  of size  $k$  and  $w \in \{0, 1\}^k$ . In the common case where  $n$  is a small constant multiple of  $k$ , we will often be somewhat loose and not explicitly distinguish between  $\mathcal{C}_k$  and  $\mathcal{C}_n$  and just refer to  $\mathcal{C}$ . Indeed, one can consider  $\mathcal{C}_k$  to effectively be a subset of  $\mathcal{C}_n$  using only some  $k$  out of  $n$  bits.

Most common Boolean classes are closed under pattern restrictions. The main notable exceptions are monotone classes such as monotone conjunctions.

We will need the following bound on the spectral norm of a pattern matrix in terms of the Fourier coefficients of  $f$ .

**Theorem 5.2.13** ([She11], Theorem 4.3). *Let  $f : \{0, 1\}^k \rightarrow \mathbb{R}$ , and let  $A$  be its  $(n, k, f)$ -pattern matrix. Let  $s = 2^n 2^k \binom{n}{k}$  be the number of entries in  $A$ . Then*

$$\|A\| = \frac{1}{s} \max_{S \subseteq [k]} \left\{ \sum_{j \in S} \widehat{f}(j) \left(\frac{k}{n}\right)^{|j|/2} \right\}.$$

### 5.3 Main theorem: correlational variance bounds via orthogonalizing distributions

Here we state our main theorem in its most general setting, that of orthogonalizing distributions. The theorem is a strong generalization of [MS22, Theorem 12], which proved such a result for the special case of DNF formulas, leveraging a powerful communication complexity result of Razborov and Sherstov [RS10]. Our main contribution lies in noting that this proof does not require the full strength of [RS10], and actually holds in a considerably general setting. Our lower bounds in terms of threshold and approximate degree follow as a result of this main theorem, since their dual characterizations furnish exactly the kinds of orthogonalizing distributions we need. The proof is a simple application of the pattern matrix method. It is worth noting that the functions in the hard family are all in fact juntas.

**Theorem 5.3.1.** *Let  $n > k$  be a multiple of  $k$ . Let  $f : \{0, 1\}^k \rightarrow \mathbb{R}$  be a Boolean function, and let  $\mu$  be a distribution that orthogonalizes  $f$  up to degree  $d - 1$ , i.e. a distribution such that  $\langle f, p \rangle_\mu = 0$  for all polynomials of degree less than  $d$ . There exists a  $(n, k, f)$ -pattern matrix family  $A$  of function-distribution pairs on  $\{0, 1\}^n$  such that  $\text{Var}(A) = \binom{n}{k}^d$ .*

*Proof.* For every  $(V, w) \subseteq V(n, k) \subseteq \{0, 1\}^k$ , define  $f_{V,w}(x) = f(x_V \oplus w)$  and

$D_{V,w}(x) = 2^{k-n} \mu(x_V = w)$ . One can verify that  $D_{V,w}$  is a valid distribution on  $\mathcal{F} = \{1, 1g^n\}$ :

$$\begin{aligned} \sum_{x \in \mathcal{F}} D_{V,w}(x) &= 2^{k-n} \sum_{x \in \mathcal{F}} \mu(x_V = w) \\ &= 2^{k-n} \sum_{z \in \mathcal{F}} \sum_{\substack{x \in \mathcal{F} \\ x_V = z}} \mu(x_V = w) \\ &= 2^{k-n} 2^{n-k} \sum_{z \in \mathcal{F}} \mu(z = w) \\ &= 1. \end{aligned}$$

Our family will be

$$A = \{(f_{V,w}, D_{V,w}) \mid (V, w) \in \mathcal{V}(n, k) \subseteq \mathcal{F} = \{1, 1g^k\}\}.$$

To analyze this, let  $g = 2^{k-n} f \mu$ , and observe that  $M(A)$  is precisely the (transpose of the)  $(n, k, g)$ -pattern matrix. For all  $j \in \mathcal{S}$ ,  $\chi_{\mathcal{S}}^j \mu = 0$ , i.e.  $\sum_x f(x) \chi_{\mathcal{S}}(x) \mu(x) = 0$ , we have that  $\widehat{g}(\mathcal{S}) = hg, \chi_{\mathcal{S}}^j = 0$ . Moreover, we know that for all  $\mathcal{S}$ ,

$$\widehat{g}(\mathcal{S}) = hg, \chi_{\mathcal{S}}^j = 2^{-k} \sum_{x \in \mathcal{F}} 2^{k-n} f(x) \mu(x) \chi_{\mathcal{S}}(x) = 2^{-n} \sum_{x \in \mathcal{F}} \mu(x) = 2^{-n}.$$

Thus by Theorem 5.2.13, we have that

$$\|M(A)\|_2 \leq \frac{1}{\sqrt{s}} 2^{-n} \left(\frac{k}{n}\right)^{d/2},$$

where  $s = 2^n 2^k (n/k)^k$ . Note that  $j \in \mathcal{A} = \{2^k (n/k)^k = 2^{-n} s\}$ . So finally by Lemma 5.2.4, we have

$$\text{Var}(A) \leq \frac{2^n}{j \in \mathcal{A}} \|M(A)\|_2^2 = \frac{2^n}{2^{-n} s} 2^{-2n} \left(\frac{k}{n}\right)^d = \left(\frac{k}{n}\right)^d.$$

□

The following corollary is immediate from Theorem 5.2.3.

**Corollary 5.3.2.** *Let  $A$  be the family defined above. Any SQ learner for  $A$  making only correlational queries of tolerance  $\tau < \frac{1}{10}$  must use  $\Omega\left(\left(\frac{n}{k}\right)^d \tau^2\right)$  queries in order to output a function with 0-1 error at most  $\frac{1}{3}$ .*



## 5.4 CSQ lower bounds for PAC learning in terms of threshold degree

Our main lower bound for PAC learning in terms of threshold degree is the following.

**Theorem 5.4.1.** *Let  $n > k$  be a multiple of  $k$ . Let  $f : \{0, 1\}^k \rightarrow \{0, 1\}$  be a function of threshold degree  $d$ . Let  $F$  denote the  $(n, k, f)$ -pattern matrix class. Any distribution-free SQ learner for  $F$  making only correlational queries of tolerance  $\tau < \frac{1}{10}$  must use  $\Omega\left(\left(\frac{n}{k}\right)^d \tau^2\right)$  queries in order to output a function with 0-1 error at most  $\frac{1}{3}$ .*

*Proof.* Letting  $\mu$  be the orthogonalizing distribution (on  $\{0, 1\}^k$ ) guaranteed by the dual characterization (Theorem 5.2.7), this theorem follows immediately by considering the family  $A$  from Theorem 5.3.1 and Corollary 5.3.2, since  $\text{Var}(A) = \left(\frac{n}{k}\right)^d$ .  $\square$

We can now prove our original Theorem 5.1.1 as stated.

*Proof of Theorem 5.1.1.* Let  $\mathcal{C}$  be a class closed under pattern restrictions (Definition 5.2.12), with threshold degree  $\Omega(d)$ . Pick a function  $f : \{0, 1\}^{n/2} \rightarrow \{0, 1\}$  in  $\mathcal{C}$  of threshold degree  $\Omega(d)$ . (Technically  $f$  lies in the version of  $\mathcal{C}$ , call it  $\mathcal{C}_{n/2}$ , on  $n/2$  bits. But as noted in Definition 5.2.12, one can effectively view  $\mathcal{C}_{n/2}$  as being part of  $\mathcal{C}$ . The threshold degree bound remains  $\Omega(d)$  asymptotically.) Now the class  $F$  constructed above (Theorem 5.4.1, with  $k = n/2$ ) must satisfy  $F \subseteq \mathcal{C}$ . The theorem follows.  $\square$

One application of this theorem, already proved as a special case in [MS22], is a  $2^{\Omega(n^{1/3})}$  SQ lower bound for learning DNF formulas on  $\{0, 1\}^n$ , proved by leveraging the well-known  $\Omega(n^{1/3})$  Minsky–Papert threshold degree bound for DNFs [MP69]. This bound matches up to logarithmic factors in the exponent the algorithm of [KS04] for this problem.

Another application is to the circuit class  $AC^0$ , for which the following essentially-optimal threshold degree bound was recently shown.

**Theorem 5.4.2** ([SW19]). *For any constant  $\delta > 0$ , there exists an  $AC^0$  circuit on  $n$  bits with threshold degree  $\Omega(n^{1-\delta})$ .*

By Theorem 5.4.1, we obtain as a consequence a CSQ lower bound of  $2^{\Omega(n^{1-\delta})}$  (for any constant  $\delta > 0$ ) for learning  $AC^0$ . This is essentially the strongest possible lower bound for distribution-free CSQ learning of  $AC^0$ .

Another important application is to intersections of two halfspaces on  $f \in \{0, 1\}^n$ , for which Sherstov [She13] showed an optimal threshold degree bound of  $\Omega(n)$ . This yields a CSQ lower bound of  $2^{\Omega(n)}$  for distribution-free learning of this class, which is again the strongest possible and improves considerably on the bound of [KS07] of  $2^{\Omega(\sqrt{n})}$  for intersections of  $\sqrt{n}$  halfspaces.

Needless to say, further applications hold for all classes with known threshold degree bounds, including symmetric functions, decision trees, and functions with known lower bounds on sensitivity [NS94, BDW02, Hua19, KSP20].

It is also worth noting that if we treat  $k$  as a free parameter, then Theorem 5.4.1 yields hard families of  $k$ -juntas. To our knowledge, the class of parities on  $k$  bits, i.e.  $f_{\chi_S} : \{0, 1\}^k \rightarrow \{0, 1\}$ , was the only class for which an optimal, exponential SQ lower bound of  $\sum_{i=0}^k \binom{n}{i} = 2^n$  was known (under the uniform distribution), matching the algorithm of [MOS03] (itself only a polynomial improvement over brute force). We see that by picking  $f : \{0, 1\}^k \rightarrow \{0, 1\}$  to be any function of threshold degree  $\Omega(k)$ , such as an intersection of two halfspaces [She13], then we can construct a family  $\mathcal{A}$  of  $k$ -juntas with a nearly-optimal distribution-free CSQ lower bound of  $2^{\Omega(\frac{n}{k})}$ .

## 5.5 CSQ lower bounds for agnostic learning in terms of approximate degree

Our lower bound for agnostic learning is obtained by exploiting the dual characterization of approximate degree (Theorem 5.2.9), which turns out to be precisely suited to this task. This is done as follows. Let  $f$  be a function with high  $\frac{1}{2}$ -approximate degree. First, since the dual function  $h$  is accompanied by an orthogonalizing distribution  $\mu$ , we can use Theorem 5.3.1 to generate a family from  $h$  that is hard to learn. Next, since the dual function  $h$  and the original  $f$  are well-correlated, specifically  $\langle hf, h \rangle_{\mu} > \frac{1}{2}$ , we can use an agnostic learner for a class generated from  $f$  to yield a (weak) PAC learner for the family generated from  $h$ . Since the latter problem is hard, the original problem is as well.

**Lemma 5.5.1.** *Let  $n > k$  be a multiple of  $k$ . Let  $f : \{0, 1\}^k \rightarrow \{0, 1\}$  be such that  $\widetilde{\deg}_{\epsilon}(f) = d$ , and let  $h : \{0, 1\}^k \rightarrow \{0, 1\}$  and the orthogonalizing distribution  $\mu$  on  $\{0, 1\}^k$  be as given by its dual characterization (Theorem 5.2.9). There exists a  $(n, k, h)$ -pattern matrix family  $A$  of function-distribution pairs on  $\{0, 1\}^n$  such that any SQ learner for  $A$  making only correlational queries of tolerance  $\tau < \frac{1}{10}$  must use  $\Omega\left(\left(\frac{n}{k}\right)^d \tau^2\right)$  queries in order to output a function with 0-1 error at most  $\frac{1}{3}$ .*

*Proof.* This is immediate from Theorem 5.3.1 and Corollary 5.3.2. □

This lets us prove an agnostic learning hardness result by reducing the problem of PAC learning  $A$  to agnostically learning a similar pattern matrix family generated from  $f$ .

**Lemma 5.5.2.** *Let  $f : \{0, 1\}^k \rightarrow \{0, 1\}$  be such that  $\widetilde{\deg}_{1-\alpha}(f) = d$ . Let  $F$  denote the  $(n, k, f)$ -pattern matrix class. Then a distribution-free agnostic learner for  $F$  capable of achieving 0-1 error  $\text{opt} + \epsilon$  yields a PAC learner capable of achieving 0-1 error  $\alpha/2 + \epsilon$  for the family  $A$  in Lemma 5.5.1.*

*Proof.* For brevity let  $D = D_{V,w}$  and let  $f^\theta, h^\theta = f_{V,w}, h_{V,w}$ . It is not hard to see that  $hf^\theta, h^\theta i_D = hf, hi_\mu$ :

$$\begin{aligned}
hf^\theta, h^\theta i_D &= \mathbb{E}_x [f^\theta(x)h^\theta(x)] \\
&= 2^{k-n} \sum_{x \in \mathcal{Z}^{2^f-1}, 1, 1g^n} f(x_V-w)h(x_V-w)\mu(x_V-w) \\
&= 2^{k-n} \sum_{z \in \mathcal{Z}^{2^f-1}, 1, 1g^k} \sum_{\substack{x \in \mathcal{Z}^{2^f-1}, 1, 1g^n: \\ x_V=z}} f(x_V-w)h(x_V-w)\mu(x_V-w) \\
&= 2^{k-n} 2^{n-k} \sum_{z \in \mathcal{Z}^{2^f-1}, 1, 1g^k} f(z-w)h(z-w)\mu(z-w) \\
&= \sum_{z \in \mathcal{Z}^{2^f-1}, 1, 1g^k} f(z)h(z)\mu(z) \\
&= hf, hi_\mu.
\end{aligned}$$

So  $hf^\theta, h^\theta i_D = hf, hi_\mu > 1 - \alpha$ , by Theorem 5.2.9. This means  $\mathbb{P}_{x \sim D}[f^\theta(x) \neq h^\theta(x)] = (1 - hf^\theta, h^\theta i_D)/2 < \alpha/2$ . In other words, each function-distribution pair  $(h^\theta, D) \in \mathcal{A}$  has  $\text{opt} = \inf_{f \in \mathcal{F}} \mathbb{P}_{x \sim D}[f^\theta(x) \neq h^\theta(x)] < \alpha/2$  with respect to  $F$ . So if we could agnostically learn  $F$  with error  $\text{opt} + \epsilon$  in a distribution-free way, then we could PAC-learn  $\mathcal{A}$  with error  $\alpha/2 + \epsilon$ .  $\square$

**Theorem 5.5.3.** *Let  $n > k$  be a multiple of  $k$ . Let  $f : \mathcal{Z}^{2^f-1}, 1, 1g^k \rightarrow \mathcal{Z}^{2^f-1}, 1, 1g$  be such that  $\widetilde{\text{deg}}_{1/2}(f) = d$ . Let  $F$  denote the  $(n, k, f)$ -pattern matrix class. Any distribution-free agnostic learner for  $F$  using only correlational statistical queries of tolerance  $\tau < \frac{1}{10}$  requires at least  $\Omega\left(\left(\frac{n}{k}\right)^{d\tau^2}\right)$  queries in order to output a function with excess 0-1 error  $\frac{1}{100}$ , i.e. true 0-1 error  $\text{opt} + \frac{1}{100}$ .*

*Proof.* Suppose we had such a learner. By Lemma 5.5.2, taking  $\alpha = \frac{1}{2}$ , this would be a PAC learner for  $\mathcal{A}$  capable of achieving 0-1 error at most  $\frac{1}{4} + \frac{1}{100} < \frac{1}{3}$ . Such a learner must obey the bound in Lemma 5.5.1.  $\square$

Theorem 5.1.2 now follows in just the same way as Theorem 5.1.1.

An important application of this theorem is the problem of agnostically learning conjunctions (likewise disjunctions) in the distribution-free setting. This was one of the original problems considered in [KSS94], and has seen hardness results in various restricted settings over the years, notable among which are a  $2^{\binom{p}{n}}$  lower bound on Perceptron-based approaches (via the approximate rank) [KS10] and a super-polynomial (but not exponential) CSQ lower bound for learning monotone conjunctions under the uniform distribution [Fel12]. Since the  $\frac{1}{2}$ -approximate degree of conjunctions on  $n$  bits is  $\Theta(\binom{p}{n})$  [NS94, Pat92], we obtain a CSQ lower bound of  $2^{\binom{p}{n}}$  for this problem. This is essentially the best possible CSQ lower bound for this problem.

Another important application is the problem of agnostically learning halfspaces [KKMS08]. Halfspaces can compute Majority functions on  $n$  bits, which have an approximate degree of  $\Omega(n)$  [Pat92]. This yields a  $2^{\binom{p}{n}}$  CSQ lower bound distribution-free agnostic learning of halfspaces, which is the strongest possible bound. Plugging in known approximate degree bounds for other classes such as symmetric Boolean functions [Pat92, BDW02] yield further applications.

One can also apply this theorem to the class  $AC^0$ , even though the resulting theorem is already implied by the PAC lower bound proved via threshold degree earlier. It is interesting that this bound can also be proved directly, via the following approximate degree bound.

**Theorem 5.5.4** ([BT19]). *For any constant  $\delta > 0$ , there exists an  $AC^0$  circuit on  $n$  bits with  $\frac{1}{2}$ -approximate degree  $\Omega(n^{1-\delta})$ .*

This yields a CSQ lower bound of  $2^{\binom{p}{n^{1-\delta}}}$  (for any constant  $\delta$ ) for distribution-free agnostic learning of  $AC^0$ .

### 5.5.1 CSQ lower bounds on attribute-efficient agnostic learning of sparse halfspaces

The problem of attribute-efficient learning [Blu90, BL97] formalizes a notion of learning in the presence of irrelevant attributes, and is an important open problem in learning theory. Consider the class of  $k$ -sparse halfspaces on  $f \in \{-1, 1\}^n$ . Since it has VC dimension  $O(k \log n)$ , from a statistical point of view a sample complexity of  $\text{poly}(k \log n)$  suffices to learn it (both in the PAC and agnostic settings). An efficient learner which achieves this sample complexity would be called an attribute-efficient learner. But despite years of research, no distribution-free attribute-efficient learners are known for this basic class.

In the SQ setting, the appropriate analog of sample complexity is the tolerance. Specifically, it takes  $\Theta(1/\tau^2)$  samples to simulate a query of tolerance  $\tau$ , and this is sometimes known as the estimation complexity of an SQ algorithm using tolerance  $\tau$ . Accordingly, Feldman [Fel14] posed the following question as the problem of attribute-efficient SQ learning of sparse halfspaces: does there exist an SQ algorithm capable of learning  $k$ -sparse halfspaces on  $f \in \{-1, 1\}^n$  only using  $\text{poly}(n)$  queries of tolerance  $(k \log n)^{O(1)}$ ? Though the question was originally asked in the PAC setting, here we answer it in the negative in the distribution-free, agnostic, CSQ setting. The result follows readily from our agnostic learning lower bound in terms of approximate degree.

**Theorem 5.5.5.** *There exists a family  $F$  of  $k$ -sparse halfspaces on  $f \in \{-1, 1\}^n$  such that any distribution-free agnostic learner for  $F$  only using correlational queries of tolerance  $\tau < \frac{1}{10}$  requires at least  $\binom{n}{k} \tau^2$  queries in order to agnostically learn  $F$  up to excess error  $\frac{1}{10}$ . In particular, for any  $\tau = (k \log n)^{O(1)}$ , the lower bound remains  $\binom{n}{k} \tau^2$  asymptotically, meaning no attribute-efficient CSQ learner exists.*

*Proof.* Taking  $f : \{-1, 1\}^k \rightarrow \{-1, 1\}$  to be a halfspace with  $\frac{1}{2}$ -approximate degree  $\Omega(k)$ , such as Majority, this is a direct application of Theorem 5.5.3.  $\square$

## 5.6 Hardness of approximation in terms of approximate degree

In addition to CSQ lower bounds, [MS22] also prove hardness of approximation results for a family in terms of its correlational variance. Given a hypothesis class  $H$  and a family  $A$  of function-distribution pairs, a hardness of approximation result takes the form of a statement that in the worst case over a choice of  $(f, D) \in A$ , no hypothesis  $h \in H$  can even achieve nontrivial loss with respect to  $f(D)$ , i.e. approximate  $f$  wrt  $D$ . Naturally, this certainly implies that  $H$  cannot be used to learn  $A$ . Our bounds on correlational variance in terms of approximate degree have consequences for hardness of approximation by kernelized linear functions, establishing inherent limitations of kernel methods for learning  $A$ . In particular, they show that the polynomial kernel is essentially an optimal kernel map in the distribution-free setting. This amounts to a variant of a result already observed in [She11].

Let  $f, h : \{ -1, 1 \}^n \rightarrow \mathbb{R}$  be functions, and let  $D$  be a distribution on  $\{ -1, 1 \}^n$ . Let  $L_{f(D)}(h)$  denote the squared loss  $\mathbb{E}_{x \sim D} [(h(x) - f(x))^2] = \|h - f\|_D^2$ . Fix any embedding (or kernel feature map)  $\psi : \{ -1, 1 \}^n \rightarrow [ -1, 1 ]^N$  for any  $N, B > 0$ , and define the kernelized linear class

$$H_{\psi, B} = \{ \langle w, \psi(x) \rangle : \|w\|_2 \leq B \}.$$

Then the following bound on approximation by  $H_{\psi, B}$  holds.

**Theorem 5.6.1** ([MS22]). *Let  $A$  be a family of function-distribution pairs on  $\{ -1, 1 \}^n$ .*

*Then*

$$\max_{(f, D) \in A} \min_{h \in H_{\psi, B}} L_{f(D)}(h) \geq \mathbb{E}_{(f, D) \in A} \min_{h \in H_{\psi, B}} L_{f(D)}(h) > \frac{1}{2} B^{\rho} \sqrt{\text{Var}(A)}.$$

A similar but more involved result can also be proved for approximation by depth-two neural networks; see [MS22, Theorem 3].

We obtain the following consequence.

**Theorem 5.6.2.** *Let  $\mathcal{C}$  be a Boolean function class closed under pattern restrictions, with  $(1 - \alpha)$ -approximate degree  $d$ , where  $\alpha$  is a sufficiently small constant ( $\alpha = 1/16$  suffices). Let  $H_{\psi,B}$  be the linear hypothesis class defined above. Then there exists a family  $A$  of function-distribution pairs, with the functions lying in  $\mathcal{C}$ , such that*

$$\max_{(f,D) \in A} \min_{h \in H_{\psi,B}} L_{f(D)}(h) - \mathbb{E}_{(f,D) \in A} \min_{h \in H_{\psi,B}} L_{f(D)}(h) > \frac{1}{8} \frac{B^{\rho}}{N^2} \quad (d).$$

*Proof.* Let  $f : \{0,1\}^{n/2} \rightarrow \{0,1\}$  be a function in  $\mathcal{C}$  with  $(1 - \alpha)$ -approximate degree  $d$ . Let  $g : \{0,1\}^{n/2} \rightarrow \{0,1\}$  and  $\mu$  on  $\{0,1\}^{n/2}$  be the accompanying dual function and distribution respectively, so that  $\langle hf, gi \rangle_{\mu} = 1 - \alpha$ .

Let  $A$  denote the  $(n, n/2, f)$ -pattern matrix family of function-distribution pairs arising from  $f$  and  $\mu$ , i.e. pairs of the form  $(f_{V,w}, D_{V,w})$ , where  $f_{V,w}(x) = f(x_V \parallel w)$  and  $D_{V,w}(x) = 2^{-n/2} \mu(x_V \parallel w)$ . Similarly let  $B$  denote the  $(n, n/2, g)$ -pattern matrix arising from  $g$  and  $\mu$ . Recall that  $\text{Var}(B) = 2^{-d}$ , since  $\mu$  orthogonalizes  $g$  up to degree  $d$ .

Fix any  $V, w$ , and let  $f^0 = f_{V,w}, g^0 = g_{V,w}, D^0 = D_{V,w}$  for brevity. It is easily calculated that  $\langle hf^0, g^0 i \rangle_{D^0} = \langle hf, gi \rangle_{\mu} = 1 - \alpha$ , so that  $\langle kf^0 - g^0 k_{D^0}^2, g^0 i \rangle_{D^0} = 2 - 2\langle hf^0, g^0 i \rangle_{D^0} = 2\alpha$ . Now fix any  $h \in H_{\psi,B}$ . We will argue that if  $h$  can approximate  $f^0$  under  $D^0$ , then it can also approximate  $g^0$ , simply because  $f^0$  and  $g^0$  are close under  $D^0$ . Indeed,

$$\begin{aligned} L_{g^0(D^0)}(h) &= \langle hg^0, h k_{D^0}^2 \rangle_{D^0} \\ &= \langle kf^0 - h k_{D^0}^2 + 2kf^0 - g^0 k_{D^0}^2, hg^0 \rangle_{D^0} \\ &= 2L_{f^0(D^0)}(h) + 4\alpha, \end{aligned}$$

i.e.,

$$L_{f^0(D^0)}(h) \leq \frac{1}{2} L_{g^0(D^0)}(h) - 2\alpha.$$

Thus for every  $(f^0, D^0) \in A$ , there is a corresponding  $(g^0, D^0) \in B$  such that

$$\min_{h \in H_{\psi,B}} L_{f^0(D^0)}(h) \leq \frac{1}{2} \min_{h \in H_{\psi,B}} L_{g^0(D^0)}(h) - 2\alpha,$$



and in fact this correspondence is one-to-one. So, taking the average,

$$\begin{aligned} \mathbb{E}_{(f^\theta, D^\theta)} \min_{h \in \mathcal{H}_{\psi, B}} L_{f^\theta(D^\theta)}(h) &= \frac{1}{2} \mathbb{E}_{(g^\theta, D^\theta)} \min_{h \in \mathcal{H}_{\psi, B}} L_{g^\theta(D^\theta)}(h) - 2\alpha \\ &> \frac{1}{4} \frac{B^2}{2} \frac{1}{N} \sqrt{\text{Var}(B)} - 2\alpha \\ &= \frac{1}{8} \frac{B^2}{2} N^{-d} \end{aligned}$$

for  $\alpha = 1/16$ . □

Thus we see that for any class  $\mathcal{C}$  with approximate degree that is even super-logarithmic, i.e.  $\omega(\log n)$ , so that  $2^{-\omega(d)}$  is negligible, i.e.  $n^{-\omega(1)}$ , then no kernelized linear function (with  $B, N$  being  $\text{poly}(n)$ ) can even weakly approximate  $\mathcal{C}$ . If  $B$  is fixed, then the dimension  $N$  of the embedding must be taken to be  $2^{\omega(d)}$  for approximation to even be possible. This is essentially tight, since  $d$  is of course the approximate degree, and so approximation is certainly possible using the polynomial kernel up to degree  $d$  (so that  $N = n^{O(d)}$ ). Thus in this sense the polynomial kernel is nearly optimal for distribution-free approximation.

## 5.7 A hard class of $p$ -concepts under the uniform distribution

The correlational variance framework extends straightforwardly from Boolean functions to  $p$ -concepts. A  $p$ -concept  $c$  on  $\mathcal{F} = \{1, 1\}^n$  is a rule that to every  $x \in \mathcal{F} = \{1, 1\}^n$  assigns a random label  $y \in \mathcal{F} = \{1, 1\}^n$  such that  $\mathbb{E}[y|x] = c(x)$ . Given a distribution  $D$  on  $\mathcal{F} = \{1, 1\}^n$ , we will use  $c(D)$  to denote the distribution on  $\mathcal{F} = \{1, 1\}^n \times \mathcal{F} = \{1, 1\}^n$  of labeled points  $(x, y)$  such that  $x \in D$  and  $\mathbb{E}[y|x] = c(x)$ . The 0-1 loss of a function  $h : \mathcal{F} = \{1, 1\}^n \rightarrow \mathcal{F} = \{1, 1\}^n$  with respect to a  $p$ -concept  $c$  over  $D$  is given by  $\mathbb{P}_{(x,y) \sim c(D)}[h(x) \neq y]$ . Since the labels are random, in general no Boolean function can achieve 0 loss. The Boolean function that best fits  $c$  is given by  $\text{sign}(c)$ , and it achieves loss

$$\mathbb{P}_{(x,y) \sim c(D)}[\text{sign}(c(x)) \neq y] = \frac{1}{2} \frac{1}{2} \mathbb{E}_{(x,y) \sim c(D)}[\text{sign}(c(x))y] = \frac{1}{2} \frac{1}{2} \langle c, \text{sign}(c) \rangle_D = \frac{1}{2} \frac{1}{2} \|c\|_{1,D},$$

where  $kck_{1,D} = \mathbb{E}_x \mathbb{E}_D [j_c(x)]$ . Thus the 1-norm  $kck_{1,D}$  characterizes the best achievable advantage over random guessing.

Let  $A$  be a family of pairs  $(c, D)$ , where  $c$  is a  $p$ -concept and  $D$  is a distribution on  $\{0, 1\}^n$ . We will work in the distribution-specific setting where  $D$  is the same for all  $c$ , and is in fact just the uniform distribution  $U$ . In this case we know that CSQ is equivalent to general SQ. One can now define  $\text{Var}(A)$  exactly as before, and exactly the same proof as the one for Theorem 5.2.3 establishes the following theorem.

**Theorem 5.7.1.** *Let  $A$  be a family as above. Consider a learner given SQ access to an unknown  $(c, D) \in A$  with the goal of finding  $h$  such that the 0-1 error  $\mathbb{P}_{x \sim c(D)} [h(x) \neq y]$  is small. Then any such learner making only queries of tolerance  $\tau$  must use  $\Omega(\frac{\tau^2}{\text{Var}(A)})$  queries in order to output a function with 0-1 error better than  $\frac{1}{2} - \frac{\tau}{2}$ .*

We can now prove a theorem very similar to Theorem 5.3.1, constructing a hard family of  $p$ -concepts from a function  $f : \{0, 1\}^k \rightarrow \{0, 1\}$  and an orthogonalizing distribution  $\mu$ . In essence, instead of considering pairs of the form  $(f, \mu)$ , we consider pairs of the form  $(f\mu, U)$ . Rather than consider  $f\mu$  directly, it will be convenient to work with a scaling instead, so that we have greater control over the 1-norm. This requires an assumption on the dual distribution  $\mu$ , namely that it satisfy a pointwise upper bound of the form  $\mu(x) \leq \lambda 2^{-k}$  for all  $x \in \{0, 1\}^k$ . Notice that we may always naively take  $\lambda = 2^k$ . But in some cases it may be possible to pick  $\lambda$  smaller, which would allow our bounds to improve.

**Theorem 5.7.2.** *Let  $n > k$  be a multiple of  $k$ . Let  $f : \{0, 1\}^k \rightarrow \mathbb{R}$  be a Boolean function, and let  $\mu$  be a distribution that orthogonalizes  $f$  up to degree  $d - 1$ , i.e. a distribution such that  $\langle f, p \rangle_\mu = 0$  for all polynomials of degree less than  $d$ . Assume that  $\mu(x) \leq \lambda 2^{-k}$  for all  $x \in \{0, 1\}^k$ . Fix the underlying distribution on  $\{0, 1\}^n$  to be uniform. There exists an  $(n, k, f\mu)$ -pattern matrix family  $A$  of  $p$ -concepts on  $\{0, 1\}^n$  such that  $\text{Var}(A) \leq \frac{1}{\lambda^2} \binom{k}{n}^d$ .*

*Proof.* Let  $U$  denote the uniform distribution on  $\{f : \{1, 1\}^{\mathcal{G}^n}\}$ . For every  $(V, w) \in \mathcal{V}(n, k) \subseteq \{f : \{1, 1\}^{\mathcal{G}^k}\}$ , let  $f_{V,w}(x) = f(x_V \cup w)$ ,  $\mu_{V,w}(x) = \mu(x_V \cup w)$ , and define  $h_{V,w}(x) = \frac{2^k}{\lambda} f(x_V \cup w) \mu(x_V \cup w)$ . Since  $\mu = \lambda 2^{-k}$ ,  $h_{V,w}$  is a well-defined  $p$ -concept, and its 1-norm wrt  $U$  is given by

$$\|h_{V,w}\|_{1,U} = 2^{-n} \sum_{x \in \{1, 1\}^{\mathcal{G}^n}} \frac{2^k}{\lambda} \mu(x_V \cup w) = 2^{-n} \frac{2^k}{\lambda} 2^{n-k} = \frac{1}{\lambda}.$$

Its best-fitting Boolean function is  $\text{sign}(h_{V,w}) = f_{V,w}$ . Our family will be

$$A = \{(h_{V,w}, U) \mid (V, w) \in \mathcal{V}(n, k) \subseteq \{f : \{1, 1\}^{\mathcal{G}^k}\}\}.$$

To analyze this, let  $g = \frac{2^{k-n}}{\lambda} f \mu$ , and observe that  $M(A)$  is precisely the (transpose of the)  $(n, k, g)$ -pattern matrix. For all  $j \in S$ , since  $\langle h_{V,w}, \chi_S \rangle = 0$ , i.e.  $\sum_x f(x) \chi_S(x) \mu(x) = 0$ , we have that  $\widehat{g}(S) = \langle h_{V,w}, \chi_S \rangle = 0$ . Moreover, we know that for all  $S$ ,

$$\widehat{g}(S) = \langle h_{V,w}, \chi_S \rangle = 2^{-k} \sum_{x \in \{1, 1\}^{\mathcal{G}^k}} \frac{2^{k-n}}{\lambda} f(x) \mu(x) \chi_S(x) = \frac{2^{-n}}{\lambda} \sum_{x \in \{1, 1\}^{\mathcal{G}^k}} \mu(x) = \frac{2^{-n}}{\lambda}.$$

Thus by Theorem 5.2.13, we have that

$$\|M(A)\|_F = \frac{2^{-n}}{\lambda} \left(\frac{k}{n}\right)^{d/2},$$

where  $s = 2^n 2^k (n/k)^k$ . Note that  $|A| = 2^k (n/k)^k = 2^{-n} s$ . So finally by Lemma 5.2.4, we have

$$\text{Var}(A) = \frac{2^n}{|A|} \|M(A)\|_F^2 = \frac{2^n}{2^{-n} s} \frac{2^{-2n}}{\lambda^2} \left(\frac{k}{n}\right)^d = \frac{1}{\lambda^2} \left(\frac{k}{n}\right)^d.$$

□

We can now prove hardness of agnostic learning in terms of approximate degree by showing that an agnostic learner could learn a hard family of  $p$ -concepts constructed as above. Specifically, let  $f : \{f : \{1, 1\}^{\mathcal{G}^k}\} \subseteq \{f : \{1, 1\}^{\mathcal{G}^n}\}$  be such that  $\frac{1}{2}$ -approximate degree is  $\widetilde{\text{deg}}_{1/2}(f) = d$  (we stress that this is a function of  $k$ , and perhaps best written

$d(k)$ ). Let  $h : \mathcal{F}^{-1}, 1g^k \rightarrow \mathcal{F}^{-1}, 1g$  and the orthogonalizing distribution  $\mu$  on  $\mathcal{F}^{-1}, 1g^k$  be as given by its dual characterization (Theorem 5.2.9). Assume that  $\mu(x) \leq \lambda 2^{-k}$  for all  $x \in \mathcal{F}^{-1}, 1g^k$ . Define  $g_{V,w}(x) = \frac{2^k}{\lambda} h(x_V - w) \mu(x_V - w)$ . Consider the  $p$ -concept family

$$A = \{(g_{V,w}, U) \mid (V, w) \in \mathcal{V}(n, k) \subseteq \mathcal{F}^{-1}, 1g^k\}.$$

By the previous theorems, we know that any learner for  $A$  using queries of tolerance  $\tau$  requires  $\Omega(\frac{\tau^2}{\text{var}(A)}) = \Omega((\frac{n}{k})^d \lambda^2 \tau^2)$  queries in order to output a function with 0-1 error better than  $\frac{1}{2} - \frac{\tau}{2}$ .

**Theorem 5.7.3.** *Let  $F$  denote the  $(n, k, f)$ -pattern matrix class. Any agnostic learner for  $F$  under  $U$  using only statistical queries of tolerance  $\tau < \frac{1}{8\lambda}$  requires at least  $\Omega((\frac{n}{k})^d \lambda^2 \tau^2)$  queries in order to output a function with excess 0-1 error  $\frac{1}{8\lambda}$ , i.e. true 0-1 error  $\text{opt} + \frac{1}{8\lambda}$ .*

*Proof.* Suppose we had such a learner. Take the class of  $p$ -concepts  $A$  defined above, and consider any  $g_{V,w}$  in  $A$ . Let  $g^\theta = g_{V,w}$  for brevity. What is the optimal error  $\text{opt} = \inf_{f^\theta \in F} \mathbb{P}_x \mathbb{E}_U [f^\theta(x) \neq g^\theta(x)]$  that any function in  $F$  achieves in approximating  $g^\theta$  under  $U$ ? This is the same as  $\frac{1}{2} - \frac{1}{2} \sup_{f^\theta \in F} \langle f^\theta, g^\theta \rangle_U$ . Observe that  $f^\theta = f_{V,w}$  does quite well:

$$\begin{aligned} \langle f^\theta, g^\theta \rangle_U &= 2^{-n} \sum_{x \in \mathcal{F}^{-1}, 1g^n} f(x_V - w) g(x_V - w) \\ &= \frac{2^k}{\lambda} 2^{-n} \sum_{x \in \mathcal{F}^{-1}, 1g^n} f(x_V - w) h(x_V - w) \mu(x_V - w) \\ &= \frac{1}{\lambda} \sum_{z \in \mathcal{F}^{-1}, 1g^k} f(z) h(z) \mu(z) \\ &= \frac{1}{\lambda} \langle f, h \rangle_\mu \\ &= \frac{1}{2\lambda}, \end{aligned}$$

since  $hf, hi_\mu = 1/2$  (by the dual characterization). Thus  $\text{opt} = \frac{1}{2} - \frac{1}{4\lambda}$ . Now, if  $\epsilon, \tau = \frac{1}{8\lambda}$ , then our agnostic learner is able to achieve 0-1 error

$$\text{opt} + \epsilon = \frac{1}{2} - \frac{1}{8\lambda} + \frac{1}{8\lambda} = \frac{1}{2} - \frac{\tau}{2}$$

when learning  $A$ . Therefore it must obey the  $\Omega\left(\left(\frac{n}{k}\right)^d \lambda^2 \tau^2\right)$  lower bound on learning  $A$ .  $\square$

Overall, we see that the applicability of this result boils down to what  $\lambda$  we can actually achieve. We also know that this is actually a lower bound on uniform-distribution agnostic learning, so it can never improve on the [KKMS08] upper bound. We can trivially always take  $\lambda = 2^k$ , which corresponds to not scaling the  $p$ -concepts at all. We do, however, obtain some nontrivial results when we take  $k = \Theta(\log n)$ . For instance, for intersections of two halfspaces, we know that  $d = \Theta(k)$  [She13], so the lower bound we get would be  $\left(\frac{n}{\log n}\right)^{\Theta(\log n)}$ , i.e. quasipolynomial.

# Chapter 6: A Moment-Matching Approach to Testable Learning and a New Characterization of Rademacher Complexity

## 6.1 Introduction

In the fundamental model of agnostic learning [KSS92, Vap98], a learner tries to output the best-fitting function from a concept class  $\mathcal{C}$  with respect to an unknown labeled distribution  $D$  in the following sense: given sufficiently many labeled examples, with high probability it must produce a hypothesis with error at most  $\text{opt}(\mathcal{C}, D) + \epsilon$  over  $D$ , where  $\text{opt}(\mathcal{C}, D)$  denotes the optimal error achievable over  $D$  by any concept in  $\mathcal{C}$ . No assumptions are made on the labels.

Agnostic learning is known to be computationally intractable for even the simplest function classes without making assumptions on the marginal [KSS92, KV94a, KS09, GR09, FGRW12, Dan16a, DSS16]. There is now a substantial literature of efficient agnostic learning algorithms under various distributional assumptions, the most common being that the marginal is Gaussian or  $\text{Unif}^{\mathcal{F}} \mathbb{1}_{\mathcal{G}^d}$  (see e.g. [LMN93, BT96, KKMS08, KOS08, Kan11]). The problem of directly verifying this distributional assumption from samples, however, is often computationally infeasible (such as for  $\text{Unif}^{\mathcal{F}} \mathbb{1}_{\mathcal{G}^d}$ ) or fundamentally ill-posed (as for  $\mathcal{N}(0, I_d)$ <sup>1</sup>).

Since in the agnostic model we make no assumptions on the labels, we have no a priori estimate of  $\text{opt}(\mathcal{C}, D)$ , the error of the best-fitting classifier. Thus, a major (and often overlooked) issue with the agnostic learning model is that *it is unclear how to verify that the agnostic learner has actually succeeded*. Note that while we can

---

<sup>1</sup>To see why this is the case even when  $d = 1$ , fix any finite sample size  $m$ , and consider a (random) discrete distribution  $\hat{D}$  that is uniform on  $\omega(m^2)$  points drawn from  $\mathcal{N}(0, 1)$ . This distribution has TV distance 1 from  $\mathcal{N}(0, 1)$  (since it is discrete), yet a sample of size  $m$  drawn from  $\hat{D}$  will have  $o(1)$  TV distance from a sample of size  $m$  drawn directly from  $\mathcal{N}(0, 1)$ .

estimate the true error of the output hypothesis on a hold-out set (a.k.a. validation), we do not know its relationship to  $\text{opt}(\mathcal{C}, D)$ .

With this motivation in mind, very recent work of Rubinfeld and Vasilyan [RV23] introduced the elegant model of testable agnostic learning, or just testable learning for short. In this model, no assumptions are made on  $D$ , but there is a tester responsible for verifying whether the unknown marginal is suitably well-behaved. Whenever the tester accepts, the learner must succeed at producing a hypothesis with error at most  $\text{opt}(\mathcal{C}, D) + \epsilon$  (with high probability). And to ensure nontriviality, whenever the unknown marginal is indeed a certain well-behaved target marginal  $D_X$ , the tester must accept (with high probability). We say the class  $\mathcal{C}$  is testably learnable with respect to a target marginal  $D_X$  if there is a tester-learner pair meeting these conditions (see Definition 6.2.3).

In this model, [RV23] showed that halfspaces can be testably learned with respect to the standard Gaussian as well as the uniform distribution on the hypercube in time and sample complexity  $d^{\tilde{O}(1/\epsilon^4)}$ . Their algorithm involves checking that the low-degree moments of the unknown marginal are close to those of the target marginal  $D_X$ . For the case where the target marginal is the standard Gaussian, they show that this implies concentration and anticoncentration properties of the unknown marginal and further prove that any distribution that satisfies such properties admits low-degree polynomial approximators for halfspaces. This analysis, however, is catered specifically to the case of halfspaces and the Gaussian target marginal. For the case where the target marginal is the uniform distribution on the hypercube, their proof combines the construction for the Gaussian case with the “critical index” framework of [DGJ<sup>+</sup>10].

### 6.1.1 Our results

Our main algorithmic contribution is a general framework that yields efficient testable learning algorithms for broad classes of functions and distributions (both

continuous and discrete). As we discuss in more detail below, our framework departs from the focus on constructing low-degree polynomial approximations with respect to absolute loss as in [RV23], which appears hard to extend to classes beyond a single halfspace. Instead, we rely on a new connection to a stronger type of approximator — sandwiching polynomials — that arises naturally in constructing pseudorandom generators for classes of Boolean functions.

As it turns out, many interesting and well-studied concept classes admit both sandwiching approximators and ordinary low-degree polynomial approximators of essentially the same degree, even though sandwiching is a formally stronger notion. As a result, we derive testable learning algorithms for halfspaces and more generally arbitrary functions of a bounded number of halfspaces with respect to any fixed strongly logconcave distribution. For the uniform distribution on the hypercube, we obtain algorithms for halfspaces, degree-2 PTFs, and constant-depth circuits. For each of these applications, our running times and sample complexity guarantees match the best known results for ordinary agnostic learning, thus showing that testable learning can often be achieved at no additional cost (see Theorem 6.4.1 and Theorem 6.5.2 for precise statements).

In particular, for the special case of testably learning a single halfspace with respect to the Gaussian or the uniform distribution on the hypercube, our results improve the  $d^{\tilde{O}(1/\epsilon^4)}$  running time and sample complexity guarantee shown in [RV23] to  $d^{\tilde{O}(1/\epsilon^2)}$ , matching the best known (and conditionally optimal) results for ordinary agnostic learning. Moreover, our analysis extends to a broad family of distributions including strongly logconcave distributions.

We now describe our results and discuss our techniques in more detail.

**Sandwiching polynomial approximation.** Our starting point is a relationship between testable learning and a certain stronger notion of polynomial approximation that arises naturally in building pseudorandom generators for Boolean func-



tion classes. Specifically, a concept class  $\mathcal{C}$  admits *sandwiching* approximations of degree  $k$  and error  $\epsilon$  on  $D_X$  if for every function  $f \in \mathcal{C}$ , there exist two degree- $k$  polynomials  $p_l, p_u$  such that for every  $x$ ,  $p_l(x) \leq f(x) \leq p_u(x)$ , and moreover  $\mathbb{E}_{D_X}[f - p_l], \mathbb{E}_{D_X}[p_u - f] \leq \epsilon$ . Observe that this is a stronger requirement than the existence of approximating polynomials for  $\mathcal{C}$  with respect to absolute loss, which only requires that for every  $f$ , there be a degree  $k$  polynomial  $p$  such that  $\mathbb{E}_{D_X}[|p - f|] \leq \epsilon$ .

The main theorem underlying our framework shows that unlike the existence of polynomial approximators with respect to absolute loss, the existence of sandwiching approximations universally translates into testable learning algorithms:

**Theorem 6.1.1** (Testable learning using approximate moment matching; see Theorem 6.4.5). *Let  $D_X$  be a distribution on  $X$ , and let  $\mathcal{C}$  be a concept class mapping  $X$  to  $\mathbb{R}$ . Let  $k \in \mathbb{N}, \delta \in \mathbb{R}_+$  be degree and slack parameters, and let  $\epsilon > 0$  be the error parameter. Suppose that every  $f \in \mathcal{C}$  admits degree- $k$   $\epsilon$ -sandwiching polynomials  $p_l \leq f \leq p_u$  w.r.t.  $D_X$  such that  $\|p_l\|_1, \|p_u\|_1 \leq \epsilon/\delta$ , where  $\|p_u\|_1$  (resp.  $\|p_l\|_1$ ) refers to the  $\ell_1$  norm of the coefficients of  $p_u$  (resp.  $p_l$ ). Suppose also that with high probability over a sample of size  $d^{O(k)}$ , the empirical moments of degree at most  $k$  of  $D_X$  are within  $\delta$  of their true moments. Then  $\mathcal{C}$  can be testably learned w.r.t.  $D_X$  up to excess error  $O(\epsilon)$  in sample and time complexity  $d^{O(k)}/\text{poly}(\epsilon)$ .*

Theorem 6.1.1 relies on a simple tester: verify that the empirical moments of degree at most  $k$  are close enough to the those of  $D_X$ . The correctness of the tester relies on the claim that sandwiching polynomials for  $\mathcal{C}$  under  $D_X$  are also sandwiching polynomials for  $\mathcal{C}$  under the uniform distribution  $\hat{D}$  on a large enough sample from  $D_X$ , with an additional error that scales proportional to the  $\ell_1$  norm of the coefficients of the polynomial approximators. Thus, whenever we have sandwiching approximators with appropriate bounds on the  $\ell_1$  norm of the coefficient vectors, our testable learner can simply use the now-standard degree- $k$  (absolute-loss) polynomial regression algorithm of [KKMS08].

The work of [KKMS08] showed that the existence of low-degree (not necessarily sandwiching) polynomial approximators with respect to absolute loss suffices for *ordinary* agnostic learning. Our theorem on sandwiching polynomial approximators can be thought of as the natural counterpart to their condition but for *testable* agnostic learning.

We stress that merely having ordinary polynomial approximators for  $\mathcal{C}$  with respect to absolute loss under  $D_X$  — as opposed to sandwiching polynomials — would not have sufficed for Theorem 6.1.1. This is because such approximators do not readily translate to approximators for  $\mathcal{C}$  under a different distribution  $D^\theta$  that only approximately matches low-degree moments with  $D_X$  (i.e., passes our test for  $D_X$ ). The crucial role played by sandwiching approximation is that it allows such a transfer principle: sandwiching approximators for  $\mathcal{C}$  under  $D_X$  (with sufficiently small coefficients) are also sandwiching approximators for  $\mathcal{C}$  under such a  $D^\theta$  (see Corollary 6.3.3).

Our main task now reduces to constructing sandwiching polynomials with sufficiently small coefficients. Since proofs of existence of sandwiching polynomials are sometimes nonconstructive (e.g., for constant-depth circuits over the hypercube, where the existence of such polynomials follows from LP duality [Baz09]), we require new techniques to prove bounds on the  $\ell_1$  norm of the coefficients. We make progress by crucially exploiting a form of *approximate* duality between sandwiching polynomials and moment matching.

**Moment matching and sandwiching polynomials.** The duality between fooling using moment matching and the existence of sandwiching polynomials is well-known in the setting of the Boolean hypercube [Baz09], where moment matching  $\text{Unif}^{\mathcal{F}} \perp \mathcal{G}^d$  up to degree  $k$  is equivalent to  $k$ -wise independence. We need an approximate version of this duality in order to derive a bound on the coefficients of the approximating polynomials. Moreover, we need duality to hold over continuous

domains for our applications to non-discrete settings (such as Gaussian and strongly logconcave distributions). A duality relating *exact* moment matching and sandwiching approximations over general domains was proved in [KM13].

We derive the following general duality result, which tells us that approximate moment matching fools a class  $\mathcal{C}$  over  $D_X$  iff every concept in  $\mathcal{C}$  admits a pair of sandwiching polynomials with sufficiently small coefficients. Our proof relies on establishing a strong duality result for a certain *semi-infinite* linear program using tools from general conic duality [Sha01].

**Theorem 6.1.2** (Fooling using approximate moment matching ( ) sandwiching approximation; see Theorem 6.3.2). *Let  $D_X$  be a distribution on  $X$ , and let  $\mathcal{C}$  be a concept class mapping  $X$  to  $\mathbb{R}$ . Let  $k \geq 1, \delta \geq 0$  be degree and slack parameters, and let  $\epsilon > 0$  be the error parameter. The following are equivalent:*

- (Approximate moment matching fools  $\mathcal{C}$ .) For all  $f \in \mathcal{C}$  and for all distributions  $D^0$  whose moments of degree at most  $k$  are within  $\delta$  of those of  $D_X$ , we have  $|\mathbb{E}_{D^0}[f] - \mathbb{E}_{D_X}[f]| \leq \epsilon$ .
- (Existence of sandwiching polynomials with bounded coefficients for  $\mathcal{C}$ .) For all  $f \in \mathcal{C}$ , there exist degree- $k$  polynomials  $p_l, p_u$  such that  $p_l \leq f \leq p_u$  (pointwise over  $\mathbb{R}^d$ ), and

$$\mathbb{E}_{D_X}[p_u - f] + \delta \|p_u\|_1 \leq \epsilon, \quad \mathbb{E}_{D_X}[f - p_l] + \delta \|p_l\|_1 \leq \epsilon,$$

where  $\|p_u\|_1$  (resp.  $\|p_l\|_1$ ) refers to the  $\ell_1$  norm of the coefficients of  $p_u$  (resp.  $p_l$ ).

**Applications: testably learning functions of halfspaces and more.** Combining Theorems 6.1.1 and 6.1.2, we obtain a clean framework for testable learning that reduces the task to establishing that approximate low-degree moment matching fools the target concept class over the target marginal. As our main application, we

show that any function of a constant number of halfspaces over  $\mathbb{R}^d$  can be testably learned up to excess error  $\epsilon$  in sample and time complexity  $d^{\tilde{O}(1/\epsilon^2)}$  with respect to any distribution whose directional projections are sufficiently anticoncentrated and have strictly sub-exponentially decaying tails:

**Definition 6.1.3.** We say a distribution  $D_X$  on  $\mathbb{R}^d$  is anticoncentrated and has  $\alpha$ -strictly subexponential tails if the following hold:

- (a)  *$\alpha$ -strictly subexponential tails:* For all  $\|k\| = 1$ ,  $\mathbb{P}[|hx, u| > t] \leq \exp(-Ct^{1+\alpha})$  for some constant  $C$ .
- (b) *Anticoncentration:* For all  $\|k\| = 1$  and continuous intervals  $T \subseteq \mathbb{R}$ , we have  $\mathbb{P}[|hx, u| \in T] \geq C^{\ell} |T|$  for some constant  $C^\ell$ .

**Theorem 6.1.4** (Testably learning functions of halfspaces; see Theorem 6.5.2). *Let  $\mathcal{C}$  be the class of functions of a constant number of halfspaces over  $\mathbb{R}^d$ . Let  $D_X$  be a distribution that is anticoncentrated and has  $\alpha$ -strictly subexponential tails (Definition 6.1.3). Then  $\mathcal{C}$  can be testably learned w.r.t.  $D_X$  up to excess error  $\epsilon$  using sample and time complexity  $d^{\tilde{O}(\epsilon^{-(1+\alpha)/\alpha})}$ .*

We note that even for ordinary agnostic learning, the above result is an exponential improvement in the dependence on  $\epsilon$  in the degree of the sandwiching polynomial compared to prior constructions of [KM13].

On the flip side, note that even though our framework handles  $D_X$  that come from a fairly broad family, our tester must know the low-degree moments of the particular  $D_X$  with respect to which it is expected to succeed. This is true for the approach of [RV23] as well, and it is an interesting open question whether this can be relaxed.

The class of distributions that are anticoncentrated and have strictly subexponential tails is fairly general and includes Gaussians, the uniform distribution on the unit sphere, and more generally, any strongly logconcave distribution [SW14]

(and in fact all of these examples have  $\alpha = 1$ ). This latter class includes the uniform distribution on any convex body with smooth boundary [BE85] and in particular, additive Gaussian smoothening of any convex body. Theorem 6.1.4 already matches the upper bound of [KKMS08] as well as known statistical-query (SQ) [GGK20, DKZ20a, DKPZ21] and cryptographic [DKR23, Tie23] lower bounds for ordinary agnostic learning of a single halfspace with respect to the Gaussian distribution. It also generalizes and improves the main algorithmic result of [RV23], who showed such a result for a single halfspace with time and sample complexity  $d^{\tilde{O}(1/\epsilon^4)}$ .

The key technical result underlying Theorem 6.1.4 is a proof that any distribution that approximately matches degree- $\tilde{O}(\epsilon^{-(1+\alpha)/\alpha})$  moments with a distribution  $D_X$  which is anticoncentrated and has  $\alpha$ -strictly subexponential tails fools functions of halfspaces with respect to  $D_X$  (see Theorem 6.5.6). Similar to the approach of [KM13], we rely on powerful methods arising from the classical theory of moments [KR96] and metric distances in probability [Zol84, RKSF13] to show that whenever the moments of  $D_X$  are strictly sub-exponential, moment closeness implies distribution closeness in the so-called  $\lambda$ -metric (see Section 6.5.1).

We also apply our framework to immediately obtain testable learning results with respect to the  $\text{Unif}^{\mathcal{F}} \mathbb{1}_{\mathcal{G}^d}$  in time  $d^{O(k)}$  for classes  $\mathcal{C}$  that are fooled by  $k$ -wise independence, including halfspaces [DGJ<sup>+</sup>10], degree-2 PTFs [DKN10], and constant-depth circuits [Bra10], with running time and sample complexity that matches their ordinary agnostic counterparts; see Theorem 6.4.1. Over the hypercube, the fact that approximate moment matching — i.e. almost  $k$ -wise independence — suffices to fool such classes is immediate by a result due to [AGM03].

**Moments vs anticoncentration.** Theorem 6.1.4 immediately implies that one can test anticoncentration properties of all directional marginals of a broad family of distributions by checking only the low-degree moments.

**Corollary 6.1.5** (Anticoncentration from approximate moment matching; see Corollary 6.5.7). *Fix  $\epsilon > 0$  and a distribution  $D_X$  that is anticoncentrated and has  $\alpha$ -strictly subexponential tails. Let  $D^\theta$  be any distribution whose moments of degree at most  $k = \tilde{O}(\epsilon^{-(1+\alpha)/\alpha})$  match those of  $D_X$  up to an additive slack of  $d^{-\tilde{O}(k)}$ . Then for any  $\|u\| = 1$  and any continuous interval  $T \subseteq \mathbb{R}$ ,  $\mathbb{P}_{x \sim D^\theta}[\langle x, u \rangle \in T] \leq \mathbb{P}_{x \sim D_X}[\langle x, u \rangle \in T] + \epsilon$ .*

This statement relates anticoncentration phenomena to structure in low-degree moments. In particular, any distribution that matches the first degree- $\tilde{O}(1/\epsilon^2)$  moments of a strongly logconcave distribution must have all its directional marginals anticoncentrated up to an additive error of  $\epsilon$ . In addition to being a basic result in probability, such a connection relates to verifying anticoncentration of all directional marginals from a small sample. Finding verification subroutines that extend beyond Gaussian (and the uniform distribution on the sphere) have a host of applications in algorithmic robust statistics and immediately yield efficient robust algorithms for list-decodable linear regression [KKK19, RY20a] and covariance estimation [BK21, RY20b, IK22], and robust clustering [BK20, BDH<sup>+</sup>20] of mixtures for broader families of distributions than currently known.

For the specific case of Gaussian distributions (and the uniform distribution on the unit  $d$ -dimensional sphere), such a property for the case when  $T$  is an origin centered interval was first proved in a sequence of works that introduced *certifiable anticoncentration* in the context of algorithmic robust statistics [KKK19, RY20a]. Their proofs use a polynomial approximator for the “box function” (see e.g. [KKK19, Appendix A]) and show that degree- $\tilde{O}(1/\epsilon^2)$  moments are enough to ensure  $\epsilon$ -approximate anticoncentration for origin centered intervals  $T$ . A similar argument based on approximations for the box function was used by [RV23] to show that matching degree- $\tilde{O}(1/\epsilon^4)$  moments of Gaussian implies  $\epsilon$ -approximate anticoncentration for all intervals  $T$  as above. This quartic dependence in the order of moments required appears necessary in a proof that constructs polynomial approximations for

the box function. Our argument above circumvents this bottleneck in these previous techniques and recovers the  $\epsilon$ -additive error anticoncentration from matching just the degree- $\tilde{O}(1/\epsilon^2)$  moments.

**Comparison to the algorithmic technique of [RV23].** As their main algorithmic result, Rubinfeld and Vasilyan [RV23] gave a testable learning algorithm for halfspaces that uses  $d^{\mathcal{O}(1/\epsilon^4)}$  time and samples. At a high level, the chief difference between our approach and theirs is that theirs relies on explicitly constructing polynomial approximators for any distribution  $D$  that matches low-degree moments with the target  $D_X$ , whereas ours uses a “black-box” transfer principle (Corollary 6.3.3) showing that sandwiching approximators under  $D_X$  immediately yield sandwiching approximators under  $D$ .

In more detail, their algorithm uses the fact that halfspaces admit a low-degree polynomial approximator with respect to a distribution  $D$  whenever  $D$  is anticoncentrated and has subgaussian low-degree moments. In order to verify that the empirical distribution on a large enough Gaussian sample possesses these two properties, they relate anticoncentration to low-degree moments via polynomial approximators for the box function as described above. For the case where the target marginal is the uniform distribution on the hypercube, they reuse their approximator for the box function and show that it yields a polynomial approximator for *regular* halfspaces with respect to almost  $k$ -wise independent distributions. They then utilize the “critical index” framework of [DGJ<sup>+</sup>10] to reduce the case of general halfspaces to the regular case.

Such an approach is possible for the simple setting of halfspaces over the Gaussian or the uniform distribution on the hypercube (with a suboptimal quartic dependence on  $1/\epsilon$ ), but it is already unclear how to extend it to halfspaces on non-product distributions (where Fourier methods fail) or to more expressive concept classes such as functions of halfspaces.

In contrast, an appeal to sandwiching approximation allows us to extend our testable learning results to non-anticoncentrated discrete distributions such as the uniform distribution on the hypercube, more expressive concept classes such as constant depth circuits on the hypercube and functions of halfspaces on continuous distributions, and to a broad family of distributions including all strongly logconcave distributions.

**Sample complexity of testable learning and Rademacher complexity.** One of our main contributions is a complete characterization of the sample complexity of testable learning. Similar to how VC-dimension corresponds to the sample complexity of distribution-free agnostic learning, we show that Rademacher complexity is the key quantity that controls the sample complexity of testable learning. Recall that the Rademacher complexity of a class  $\mathcal{C}$  w.r.t.  $D_X$  at sample size  $m$  is given by

$$R_m(\mathcal{C}, D_X) = \mathbb{E}_{\substack{x_i, g_{i \in [m]} \\ D_X^m}} \mathbb{E}_{\sigma} \sup_{f \in \mathcal{C}} \left| \frac{1}{m} \sum_{i \in [m]} \sigma_i f(x_i) \right|.$$

This measure plays an important role in statistical learning theory since it controls the uniform convergence of empirical losses to true losses over all  $f \in \mathcal{C}$  (see Theorem 6.2.8). To our knowledge, our result is the first natural *model-based characterization* of Rademacher complexity. We obtain precise upper and lower bounds on the sample complexity of testable learning within excess error  $\epsilon$  purely in terms of Rademacher complexity:

**Theorem 6.1.6** (Rademacher complexity characterizes testable learning, see Theorems 6.6.1 and 6.6.2). *Let  $D_X$  be a distribution on  $X$ , let  $\mathcal{C}$  be a concept class mapping  $X$  to  $\mathbb{R}$ , and let  $\epsilon > 0$  be the error parameter.*

- (Upper bound.) *Let  $m$  be such that  $R_m(\mathcal{C}, D_X) \leq \epsilon/5$ . Then  $\mathcal{C}$  can be testably learned w.r.t.  $D_X$  up to excess error  $\epsilon$  using sample complexity  $m + O(1/\epsilon^2)$ .*



- (Lower bound.) Let  $M$  be such that  $R_M(C, D_X) \leq 5\epsilon$ , and assume  $M = \Theta(1/\epsilon^2)$ . Then the sample complexity required to testably learn  $C$  w.r.t.  $D_X$  up to excess error  $\epsilon$  is at least  $\Omega(\frac{1}{\epsilon} M)$ .

This characterization yields an interesting separation between ordinary distribution-specific agnostic learning and testable learning. For the former, while uniform convergence is always a sufficient condition, it is not necessary, as witnessed by examples such as convex sets in Gaussian space [KOS08] and monotone Boolean functions [BT96] (see Section 6.6.2.1). Indeed, the sample complexity of distribution-specific agnostic learning is known to be characterized by the  $L^1(D_X)$  metric entropy rather than the Rademacher complexity [BI91]. In contrast, it is the Rademacher complexity that provides the right characterization for testable learning. Thus, testable learning is a natural supervised learning model for which bounded Rademacher complexity, and hence uniform convergence, provides a *necessary and sufficient* condition for learning. For further discussion, see Section 6.7.1.

### 6.1.2 Subsequent work and open questions

Subsequent to the initial appearance in conference proceedings of both this work [GKK23] and that of Rubinfeld and Vasilyan [RV23], there have been a few exciting developments addressing some natural questions left open by these works. One such natural question arises from the fact that the tester-learners for halfspaces in these works inherently incur an expensive  $d^{\text{poly}(1/\epsilon)}$  dependence on  $d$  and  $1/\epsilon$ , and we know that such a dependence is necessary if we seek the “strong” agnostic learning guarantee of  $\text{opt} + \epsilon$  (even in the ordinary, non-testable setting) [GGK20, DKZ20a, DKPZ21, DKR23, Tie23]. Can we obtain more efficient algorithms by relaxing the optimality guarantee we seek? The work of [GKSV23a, DKK+23] shows that the answer is yes: for the case where the target marginal is Gaussian, one can in fact design fully polynomial-time tester-learners for origin-centered halfspaces with the weaker agnostic guarantee of  $O(\text{opt})$ . Interestingly, both these works make use of the

key structural pseudorandomness result of this chapter (Theorem 6.5.6).

An even more natural and compelling question concerns the key role played by the target marginal  $D_X$  in the current framework (Definition 6.2.3). As mentioned earlier, the moment-matching tests in this work (and in fact all tester-learners discussed thus far) are inherently tailored to the particular  $D_X$  with respect to which the tester-learner is required to succeed; indeed, these tests would reject most distributions that are far from  $D_X$ , however well-behaved. Can we design tester-learners that succeed simultaneously with respect to a large class of structured distributions? This question was answered affirmatively by the very recent work of [GKSV23b], which builds on and subsumes [GKSV23a]. They show fully polynomial-time tester-learners for origin-centered halfspaces that guarantee error  $O(\text{opt})$  and accept whenever the marginal is *any* distribution that satisfies a Poincaré inequality (and some mild anticoncentration properties). This latter class of distributions includes all strongly log-concave distributions. In fact, under the Kannan–Lóvasz–Simonovits [KLS95] conjecture, it includes all log-concave distributions.

Interestingly, the natural “strong agnostic learning” version of the aforementioned result remains open. That is, can we design tester-learners for halfspaces that guarantee error  $\text{opt} + \epsilon$  (even at the cost of  $d^{\text{poly}(1/\epsilon)}$  time) while accepting a wide class  $\mathcal{T}$  of structured distributions? One approach that would suffice (by the results in this chapter) is to design a tester that checks whether the low-degree empirical moment tensor matches *any* of the low-degree moment tensors of distributions in  $\mathcal{T}$ . For  $\mathcal{T}$  being say the class of all (strongly) log-concave distributions, this is an interesting and open algorithmic question.

Many other important questions remain open. From the statistical point of view, the current characterization in terms of Rademacher complexity has a quadratic gap between the upper and lower bounds; can we provide a characterization that is tight up to constant factors? More generally, the model of testable learning is still very new, and interesting directions include designing tester-learners for richer concept

classes, as well as extending the types of guarantees that one provides testers for in the first place (error optimality guarantees are just one such type). Its connections to other areas of computational learning theory that are motivated by considerations of verifiability, such as classification with abstention [BW08] and reliable agnostic learning [KKM12], also remain to be fully explored.

### 6.1.3 Other related work

In the broader context of trustworthy or verifiable machine learning, there have been a number of works studying various different formulations of verifiability. These include: classification with abstention (see e.g. [BW08, CDM16]), in which a classifier is allowed to abstain from producing a prediction for certain points; reliable agnostic learning [KKM12], in which one seeks to control both the error as well as (say) the false positive rate (and which is closely related to a form of transductive learning [GKKM20, KK21]); and cryptographic interactive proofs for machine learning [GRSY21], in which one seeks to design an interactive proof protocol to verify the error guarantee claimed by a (potentially malicious) learner. Another very recent work that tackles issues raised by unrealistic distributional assumptions is that of [BLMT23], which proposes learning algorithms whose running times naturally scale with a measure of complexity of the underlying (unknown and arbitrary) distribution.

In pseudorandomness, the duality between fooling using bounded independence and sandwiching approximation is a fundamental tool for showing that  $k$ -wise independence fools various classes [Baz09, Bra10, DGJ<sup>+</sup>10]. Its more general statement in terms of moment matching was observed by [KM13] (see also [KKM13]), who used it to obtain low-degree sandwiching polynomials for functions of halfspaces w.r.t. logconcave distributions (their constructions do not give any insight on the  $\ell_1$  norm of the coefficients). We build on their approach for our main application, namely testably learning functions of halfspaces with respect to Gaussians, and obtain exponentially improved degree bounds in terms of  $\epsilon$  along with effective bounds on the size of the coefficients.

In statistical learning theory and nonparametric regression, one of the basic objectives is to place tight bounds on the excess risk  $L(\hat{f}) - \inf_{f \in \mathcal{C}} L(f)$  and on the generalization gap  $j\hat{L}_m(\hat{f}) - L(\hat{f})j$  of an estimator  $\hat{f}$  in various settings (see e.g. [BM97, VdG00, Tsy08]). In particular, there is a long line of work studying data-dependent bounds on these quantities in terms of measures such as the Rademacher complexity and various refinements and variants thereof [KP00, Kol01, BBL02, BM02, BBM05, Kol06]. Our sample complexity upper bound applies a simple such data-dependent bound to the testable learning setting. In terms of lower bounds, [KR21] have studied bounds on the minimal error of any ERM estimator, in the additive Gaussian noise setting, in terms of the Gaussian complexity. None of these works, however, consider a model directly comparable to testable learning.

Statistical characterizations of PAC learning have also been well-studied. In the distribution-free setting, it is very well-known that the sample complexity is characterized fully by the VC-dimension, and equivalent to uniform convergence [Vap98]. For distribution-specific agnostic setting, [BI91] obtained a characterization in terms of the  $L^1(D_X)$  metric entropy, i.e. the log covering number w.r.t. the metric  $\rho(f, g) = \mathbb{E}_{D_X}[|f - g|]$ . Work by [KL18] (see also [Vad17]) proposed a characterization of efficient agnostic learning using the so-called refutation complexity, and interpreted it as a computational analog of Rademacher complexity. Results of [SSSS10] showed that in Vapnik’s General Setting of Learning, the sample complexity is in general characterized by notions of algorithmic stability rather than uniform convergence. In modern deep learning theory, the failures of the uniform convergence paradigm in the overparameterized regime have been much studied (see e.g. [ZBH<sup>+</sup>21, NK19]); we refer the reader to [BMR21, Bel21] for surveys.

## 6.2 Preliminaries

### 6.2.1 Notation and conventions

We denote the domain by  $X$ , which for us is always either  $\mathbb{R}^d$  or  $\mathcal{F}^{-1}g$ , and labels always lie in  $\mathcal{F}^{-1}g$ . We use  $\mathcal{C}$  to denote a concept class mapping  $X$  to  $\mathcal{F}^{-1}g$ . We use  $D_X$  to denote a well-behaved distribution on  $X$  (i.e. the target marginal, such as  $\mathcal{N}(0, I_d)$  or  $\text{Unif}_{\mathcal{F}^{-1}g^d}$ ), and we use the calligraphic  $D$  to denote labeled distributions on  $X \times \mathcal{F}^{-1}g$ . We denote a size- $m$  (labeled) sample drawn from  $D$  by  $S \sim D^m$ . If  $S = \{(x_i, y_i)\}_{i \in [m]}$ , then we use  $S_X = \{x_i\}_{i \in [m]}$  to denote its “marginal”, i.e. the unlabeled sample.

Our loss function throughout will be the 0-1 loss function,  $\ell(\hat{y}, y) = \mathbb{1}[\hat{y} \neq y]$ . Given a labeled distribution  $D$ , we denote the population loss by  $L(f, D) = \mathbb{P}_{(x,y) \sim D}[f(x) \neq y]$  (or just  $L(f)$  when  $D$  is implicit), and the empirical loss on a size- $m$  sample  $S \sim D^m$  by  $\hat{L}_m(f, S) = \mathbb{P}_{(x_i, y_i) \in S}[f(x_i) \neq y_i]$  (or just  $\hat{L}_m(f)$  when  $S$  is implicit). We follow the convention of denoting empirical quantities using a hat and a subscript to denote the sample size (as in  $\hat{L}_m$ ). We use  $\text{opt}(\mathcal{C}, D)$  to denote  $\inf_{f \in \mathcal{C}} L(f, D)$ .

We follow the following conventions when working with monomials over  $X$ . For any multi-index  $I \in \mathbb{N}^d$ , let  $|I| = \sum_j I_j$  denote its degree (or sometimes order), and let  $x_I$  denote the monomial  $\prod_{j \in [d]} x_j^{I_j}$ . We use  $\mathcal{I}(k, d) = \{I \in \mathbb{N}^d \mid |I| \leq k\}$  to denote the set of multi-indices of degree at most  $k$ . For a vector  $\Delta \in \mathbb{R}_+^{|\mathcal{I}(k, d)|}$  and a degree- $k$  polynomial  $p : \mathbb{R}^d \rightarrow \mathbb{R}$  given by  $p(x) = \sum_{I \in \mathcal{I}(k, d)} \Delta_I x_I$ , we use  $\|\Delta\|_1$  to denote  $\sum_{I \in \mathcal{I}(k, d)} \Delta_I$ . This may be thought of as the  $\Delta$ -weighted  $\ell_1$  norm of the coefficients of  $p$ .

We use  $a \asymp b$  and  $a \lesssim b$  to denote equalities and inequalities up to constants. It will be convenient to state Stirling’s approximation in the following form [Rob55]: for all  $n \geq 1$ ,  $n! \asymp e^{-n} n^{n+1/2}$ . We will also make use of the double factorial  $n!! = n(n-2) \dots$ , satisfying  $n! = n!!(n-1)!!$  for even  $n$ .

Throughout this chapter, we use the term “with high probability” to mean “with probability at least 0.99” (or any other sufficiently large constant) for simplicity. In all cases, confidence may be amplified using standard repetition arguments.

## 6.2.2 Learning models

Here we formally define the learning models we work with. Let  $D_X$  be a distribution on  $X$ , where  $X$  is either  $\mathbb{R}^d$  or  $f^{-1}g^d$ , and let  $C$  be a concept class mapping  $X$  to  $f^{-1}g$ .

**Definition 6.2.1** (Distribution-specific agnostic learning). We say a learner  $A$  agnostically learns  $C$  w.r.t.  $D_X$  up to excess error  $\epsilon$  if for any  $D$  on  $X = f^{-1}g$  with marginal  $D_X$ , given sufficiently many examples drawn from  $D$ , with high probability  $A$  outputs a hypothesis  $h$  such that  $L(h) = \text{opt}(C, D) + \epsilon$ . Here, recall that  $L(f) = \mathbb{P}_{(x,y) \sim D}[f(x) \neq y]$  and  $\text{opt}(C, D) = \inf_{f \in C} L(f)$ .

We recall the standard result of [KKMS08] that shows that polynomial approximators with respect to the absolute loss yield agnostic learning algorithms.

**Theorem 6.2.2** ([KKMS08]). *Suppose that for every  $f \in C$ , there exists a degree- $k$  polynomial  $p : X \rightarrow \mathbb{R}$  such that  $\mathbb{E}_x \sim D_X[|f(x) - p(x)|] \leq \epsilon$ . Then there exists a simple agnostic learner (based on degree- $k$  polynomial regression w.r.t. the absolute loss) for learning  $C$  w.r.t.  $D_X$  up to excess error  $\epsilon$  in time and sample complexity  $d^{O(k)} / \text{poly}(\epsilon)$ .*

We now formally define testable learning.

**Definition 6.2.3** (Testable agnostic learning, [RV23]). We say a tester-learner pair  $(T, A)$  testably learns  $C$  w.r.t.  $D_X$  up to excess error  $\epsilon$  if for any distribution  $D$  on  $X = f^{-1}g$ , the following conditions are met:

- (Soundness/composability.) If  $D$  is such that the tester  $T$  accepts with high probability over a sample drawn from  $D$ , then the learner  $A$  succeeds in agnostically learning  $C$  w.r.t.  $D$ , i.e. with high probability it produces a hypothesis  $h$  such that  $L(h) = \text{opt}(C, D) + \epsilon$ .

- (Completeness.) Whenever  $D$  truly has marginal  $D_X$  on  $X$ , the tester  $T$  accepts with high probability.

Again, here “with high probability” may be taken to be “with probability at least 0.99” for simplicity, and the confidence in each step may be amplified using standard repetition arguments.

Note that as stated, it is not strictly necessary for the tester  $T$  and the learner  $A$  to work with the same sample, and the definition may be interpreted as saying “if  $T$  accepts  $D$  (w.h.p.), then  $A$  must succeed over  $D$  (w.h.p.)”. The algorithms and characterizations we give in this chapter have a stronger “data-dependent” guarantee, where both  $T$  and  $A$  operate on the same sample  $S$  drawn from  $D$ , and have the following interpretation: “if  $T$  accepts  $S$ , then  $A$  must succeed over  $S$  (as well as generalize to  $D$  w.h.p.)”.

### 6.2.3 Bounded independence and sandwiching polynomials over the hypercube

In this section, let  $U$  denote  $\text{Unif } \{0, 1\}^d$ .

**Definition 6.2.4.** We say a distribution  $D$  on  $\{0, 1\}^d$  is  $(\delta, k)$ -independent if for all  $|I| \leq k$ ,  $|J| \leq k$ ,  $J \cap I = \emptyset$ ,  $\mathbb{E}_D[x_I x_J] \leq \delta$ . When  $\delta = 0$ , we simply call  $D$  a  $k$ -wise independent distribution.

We say that a concept class  $\mathcal{C}$  is  $\epsilon$ -fooled by  $(\delta, k)$ -independence (resp.  $k$ -wise independence) if for every  $f \in \mathcal{C}$  and any  $(\delta, k)$ -independent (resp.  $k$ -wise independent)  $D$ ,  $|\mathbb{E}_D[f] - \mathbb{E}_U[f]| \leq \epsilon$ .

Notice that saying  $D$  is  $(\delta, k)$ -independent is exactly equivalent to saying that the moments of degree at most  $k$  of  $D$  are within  $\delta$  of those of  $U$  (which, of course, are all 0).

We now recall a fundamental duality result in pseudorandomness, which states that bounded independence fools a class  $\mathcal{C}$  iff it admits sandwiching polynomials w.r.t.  $U$ .

**Theorem 6.2.5** ([Baz09, Thm A.1]). *Let  $f : \mathcal{F} \rightarrow \mathbb{R}$  be a function, let  $\epsilon > 0$  be the error parameter, and let  $k \geq 2, \delta > 0$  be the degree and slack parameters. The following are equivalent:*

- (a)  $(\delta, k)$ -independence fools  $f$ .
- (b) There exist degree- $k$  polynomials  $p_l, p_u$  such that  $p_l \leq f \leq p_u$  (pointwise over  $\mathcal{F}$ ), and

$$\mathbb{E}_U[p_u - f] + \delta k p_u k_1 \leq \epsilon, \quad \mathbb{E}_U[f - p_l] + \delta k p_l k_1 \leq \epsilon,$$

where for a polynomial  $p(x) = \sum_I p_I x_I$  we use  $k p k_1$  to denote  $\sum_{I \neq \emptyset} |p_I|$ , i.e. the  $\ell_1$  norm of its (nonconstant) coefficients.

The following theorem, showing that a  $(\delta, k)$ -independent distribution is statistically close to being  $k$ -wise independent, will also be useful to us.

**Theorem 6.2.6** ([AGM03, Thm 2.1]). *Let  $D$  be a  $(\delta, k)$ -independent distribution on  $\mathcal{F} \rightarrow \mathbb{R}$ . Then there exists a  $k$ -wise independent distribution  $D^\theta$  that has TV distance at most  $\delta d^k$  from  $D$ .*

### 6.2.4 Rademacher complexity

The Rademacher complexity is one of the most well-studied measures of the complexity of a function class in statistical learning theory, and may be intuitively thought of as measuring the ability of a function class to fit a randomly-labeled sample. The following definitions and theorems are now standard in the literature (see e.g. [BBL03, BM02, Bar14, BMR21] and references therein).



**Definition 6.2.7** (Rademacher complexity). Consider a sample of  $m$  points  $S_X = \{x_i, y_i\}_{i \in [m]} \in D_X^m$ . The empirical Rademacher complexity of the class  $\mathcal{C}$  w.r.t. this sample is defined to be

$$\widehat{R}_m(\mathcal{C}, S_X) = \mathbb{E}_{\sigma} \sup_{f \in \mathcal{C}} \left| \frac{1}{m} \sum_{i \in [m]} \sigma_i f(x_i) \right|. \quad (6.2.1)$$

Note that this is a random variable depending on  $S_X$ . The (expected) Rademacher complexity of  $\mathcal{C}$  w.r.t.  $D_X$  at sample size  $m$  is defined to be

$$R_m(\mathcal{C}, D_X) = \mathbb{E}_{S_X \sim D_X^m} \widehat{R}_m(\mathcal{C}, S_X). \quad (6.2.2)$$

Sometimes we simply say  $R_m(\mathcal{C})$  (resp.  $\widehat{R}_m(\mathcal{C})$ ) when  $D_X$  (resp.  $S_X$ ) is clear from context.

The next theorem states that the Rademacher complexity of a class tightly controls uniform convergence, i.e. bounds on the quantity  $\sup_{f \in \mathcal{C}} |L(f) - \widehat{L}_m(f)|$ , where  $L$  and  $\widehat{L}_m$  are the population and empirical loss functionals. The upper bound here follows from a so-called symmetrization argument, while the lower bound follows from a desymmetrization argument. In our statement, we specialize to the case of the 0-1 loss.<sup>2</sup>

**Theorem 6.2.8** (see e.g. [Bar14]). *Let  $\mathcal{C}$  be a class of functions mapping  $X$  to  $\{0, 1\}$ . Let  $D$  be a distribution on  $X \times Y$  with marginal  $D_X$  on  $X$ , and let  $S = \{(x_i, y_i)\}_{i \in [m]}$  be a random sample of size  $m$  drawn from  $D$ . For any  $f \in \mathcal{C}$ , let  $L(f) = \mathbb{P}_{(x,y) \sim D}[f(x) \neq y]$ , and let  $\widehat{L}_m(f) = \mathbb{P}_{(x_i, y_i) \sim S}[f(x_i) \neq y_i]$ . Then with probability  $1 - \delta$  over the draw of  $S$ , we have*

$$\frac{1}{4} R_m(\mathcal{C}) \leq \mathbb{P} \left( \sup_{f \in \mathcal{C}} |L(f) - \widehat{L}_m(f)| \leq R_m(\mathcal{C}) + \Theta \left( \sqrt{\frac{\log(1/\delta)}{m}} \right) \right).$$

<sup>2</sup>For general loss functions  $\ell$ , one would define a "loss class"  $\ell \in \mathcal{C} = \{f(x, y) \mid \ell(f(x), y) \mid f \in \mathcal{C}\}$  and state such a result in terms of  $R_m(\ell \in \mathcal{C})$ . In the case of the 0-1 loss function, it is known that  $R_m(\ell \in \mathcal{C}) = \frac{1}{2} R_m(\mathcal{C})$ .

The following useful facts characterize the concentration of the quantities defining the Rademacher complexity and follow by standard applications of McDiarmid's inequality. Assume that the range of  $\mathcal{C}$  is bounded in  $[-1, 1]$ .

The first fact is that the empirical Rademacher complexity  $\widehat{R}_m(\mathcal{C})$  concentrates tightly around the expected Rademacher complexity  $R_m(\mathcal{C})$ . Formally, with probability at least  $1 - \delta$  over a sample  $S_X = \{x_i\}_{i \in [m]} \in D_X^m$  of size  $m$ , we have

$$\left| R_m(\mathcal{C}, D_X) - \widehat{R}_m(\mathcal{C}, S_X) \right| = O\left(\sqrt{\frac{\log(1/\delta)}{m}}\right). \quad (6.2.3)$$

The second is that for any fixed sample  $S$  of size  $m$ , the random variable  $\sup_{f \in \mathcal{C}} \frac{1}{m} \sum_i \sigma_i f(x_i)$  concentrates tightly around its expectation,  $\widehat{R}_m(\mathcal{C})$ . Formally, with probability at least  $1 - \delta$  over  $\sigma \in \{-1, 1\}^m$  we have

$$\left| \widehat{R}_m(\mathcal{C}, S_X) - \sup_{f \in \mathcal{C}} \frac{1}{m} \sum_i \sigma_i f(x_i) \right| = O\left(\sqrt{\frac{\log(1/\delta)}{m}}\right). \quad (6.2.4)$$

Combining Theorem 6.2.8 with Eq. (6.2.3), we actually have the following data-dependent generalization guarantee for a sample  $S$  in terms of  $\widehat{R}_m(\mathcal{C}, S_X)$  itself: with probability at least  $1 - \delta$  over the draw of  $S$ , for every  $f \in \mathcal{C}$ ,

$$\left| L(f) - \widehat{L}_m(f) \right| \leq \widehat{R}_m(\mathcal{C}) + O\left(\sqrt{\frac{\log(1/\delta)}{m}}\right). \quad (6.2.5)$$

### 6.3 Duality

In this section we state the duality between fooling using approximate moment matching and sandwiching polynomials. This is a generalization of duality over the hypercube, Theorem 6.2.5, to continuous domains and more general distributions. Note that a version of duality over  $\mathbb{R}^d$ , albeit only for exact moment matching, was stated in [KM13, Lemma 3.3].

**Definition 6.3.1** (Approximate moment matching). Let  $k \geq \mathbb{N}$  be a degree parameter, and let  $\Delta \geq \mathbb{R}_+^{j \in (k,d)^j}$  be a slack parameter, satisfying  $\Delta_0 := \Delta_{(0,\dots,0)} = 0$  and

$\Delta_I > 0$  for all other  $I \geq l(k, d)$ . We say that two distributions  $D, D^\theta$  on  $X$  match moments of degree (or order) at most  $k$  up to slack  $\Delta$  if  $|E_D[x_I] - E_{D^\theta}[x_I]| \leq \Delta_I$  for all  $I \geq l(k, d)$ .

The reason for allowing the slack  $\Delta_I$  to depend on  $I$  is that in general we expect the scale of the moments to vary widely with  $I$  (as with the Gaussian, for example). The empty index  $I_0 = 0 = (0, \dots, 0)$  plays a special role, since  $x_{I_0} = 1$  and  $E_D[1] = 1$  for any valid distribution, meaning we may assume  $\Delta_0 = 0$  without loss of generality.

We can now state the main theorem. We prove this theorem using conic LP duality [Sha01], taking care to establish strong duality, but the essential argument is similar to Bazzi's proof of Theorem 6.2.5.

**Theorem 6.3.2** (Duality). *Let  $k \geq \mathbb{N}, \Delta \geq \mathbb{R}_+^{l(k,d)}$  be the degree and slack parameters, as in Definition 6.3.1. Let  $f : X \rightarrow \mathbb{R}$  be a function, and let  $D$  be a distribution on  $X$ . The following are equivalent:*

- (a) (Approximate moment matching fools  $f$  w.r.t.  $D$ .) For any distribution  $D^\theta$  whose moments up to order  $k$  match those of  $D$  up to  $\Delta$ , we have  $|E_{D^\theta}[f] - E_D[f]| \leq \epsilon$ .
- (b) (Existence of sandwiching polynomials with bounded coefficients for  $f$  w.r.t.  $D$ .) There exist degree- $k$  polynomials  $p_l, p_u$  such that  $p_l \leq f \leq p_u$  (pointwise over  $\mathbb{R}^d$ ), and

$$E_D[p_u - f] + h\Delta, |p_u| \leq \epsilon, \quad E_D[f - p_l] + h\Delta, |p_l| \leq \epsilon.$$

(Recall that for a degree- $k$  polynomial  $p(x) = \sum_I p_I x_I$ , we use  $h\Delta, |p_I|$  to denote  $\sum_{I \geq l(k,d)} |p_I| \Delta_I$ .)

*Proof.* Let  $\sigma_I = E_D[x_I]$ . Let  $P_d$  be the set of all Borel probability measures on  $\mathbb{R}^d$ . Consider the following semi-infinite linear program, which seeks to maximize  $E_{D^\theta}[f]$

over all probability distributions  $D^\theta$  on  $\mathbb{R}^d$  that approximately match moments with  $D$ :

$$\sup_{D^\theta \in \mathcal{P}_d} \mathbb{E}_{D^\theta}[f] \quad (6.3.1)$$

$$\text{subject to } \sigma_I + \Delta_I \leq \mathbb{E}_{D^\theta}[x_I] \quad \sigma_I + \Delta_I \in \mathcal{I}(k, d) \quad (6.3.2)$$

The case of  $I = (0, \dots, 0)$  is special: here  $\sigma_0 = \mathbb{E}_D[1] = 1$  and  $\Delta_0 = 0$ , so the corresponding constraint becomes simply  $\mathbb{E}_{D^\theta}[1] = 1$ , which is equivalent to requiring that  $D^\theta$  be a valid probability measure.

The dual LP turns out to be equivalent to the following, with variable  $\beta \in \mathbb{R}^{\mathcal{I}(k, d)}$ :

$$\inf_{\beta \in \mathbb{R}^{\mathcal{I}(k, d) + 1}} \sum_{I \in \mathcal{I}(k, d)} \beta_I \sigma_I + \sum_{I \in \mathcal{I}(k, d)} j \beta_I j \Delta_I \quad (6.3.3)$$

$$\text{subject to } \sum_{I \in \mathcal{I}(k, d)} \beta_I x_I \leq f(x) \quad \forall x \in \mathbb{R}^d \quad (6.3.4)$$

Notice that the primal LP (Eq. (6.3.1)) is feasible (indeed, by  $D^\theta = D$ ), and moreover, we claim that strong duality holds. Accepting this for a moment, denote the common optimum of Eqs. (6.3.1) and (6.3.3) by  $\gamma$ . The claim that approximate moment matching fools  $f$  (in a one-sided fashion) w.r.t.  $D$  is the same as asserting  $\gamma \leq \mathbb{E}_D[f] + \epsilon$ . Take  $\beta$  to be an optimal solution to the dual, and let  $p_u(x) = \sum_{I \in \mathcal{I}(k, d)} \beta_I x_I$ . (In fact, this correspondence between degree- $k$  polynomials and their coefficient vectors allows us to equivalently view the dual as optimizing over such polynomials instead of their coefficients.) The dual then tells us that  $p_u \leq f$  pointwise, and

$$\gamma = \sum_{I \in \mathcal{I}(k, d)} \beta_I \sigma_I + \sum_{I \in \mathcal{I}(k, d)} j \beta_I j \Delta_I = \mathbb{E}_D[p_u] + \langle \Delta, j \beta j \rangle \leq \mathbb{E}_D[f] + \epsilon,$$

establishing the existence of the upper sandwiching polynomial. To obtain the lower sandwiching polynomial, we replace the objective of the primal with  $\mathbb{E}_{D^\theta}[f]$  and repeat the same argument, this time using the fact that the common optimum  $\gamma^\theta$  satisfies  $\gamma^\theta \geq \mathbb{E}_D[f] + \epsilon$  (i.e., effectively replacing  $f$  with  $-f$  throughout). This establishes the desired equivalence if we accept strong duality.

Formally, it remains to properly justify that the primal LP (Eq. (6.3.1)) is well-posed, that Eq. (6.3.3) is indeed the dual of Eq. (6.3.1), and that strong duality holds. We do so in Section D.1 using results from general conic LP duality [Sha01].  $\square$

To see how Theorem 6.1.2 may be recovered from this, simply set  $\Delta$  to be  $\delta$  in every coordinate.

For the purposes of testing using moment matching, one can only ever hope to check that the unknown marginal ( $D^\theta$ , say) *approximately* matches moments with the target marginal. Approximate duality — and specifically the appearance of the quantities  $h\Delta, jp_u j$  and  $h\Delta, jp_l j$  — turns out to be precisely what we need to guarantee sandwiching polynomials even w.r.t. such a  $D^\theta$ . This is our black-box transfer principle.

**Corollary 6.3.3** (Transfer principle). *Let  $f, D, k, \Delta, \epsilon$  satisfy the conditions of Theorem 6.3.2, and let  $p_l, p_u$  be the resulting sandwiching polynomials for  $f$  w.r.t.  $D$ . Consider any particular  $D^\theta$  whose moments up to order  $k$  match those of  $D$  up to  $\Delta$ . Then  $p_l, p_u$  are sandwiching polynomials for  $f$  w.r.t.  $D^\theta$  as well, satisfying*

$$\mathbb{E}_{D^\theta}[p_u - f] \leq 2\epsilon, \quad \mathbb{E}_{D^\theta}[f - p_l] \leq 2\epsilon.$$

*Proof.* By the first part of Theorem 6.3.2, we know  $j\mathbb{E}_D[f] - \mathbb{E}_{D^\theta}[f]j \leq \epsilon$ . Thus

$$\left| \mathbb{E}_D[p_u - f] - \mathbb{E}_{D^\theta}[p_u - f] \right| = \left| \mathbb{E}_D[p_u] - \mathbb{E}_{D^\theta}[p_u] \right| + \left| \mathbb{E}_D[f] - \mathbb{E}_{D^\theta}[f] \right| \tag{6.3.5}$$

$$h\Delta, jp_u j + \epsilon. \tag{6.3.6}$$

Applying the second part of Theorem 6.3.2, this means

$$\mathbb{E}_{D^\theta}[p_u - f] \leq \mathbb{E}_D[p_u - f] + h\Delta, jp_u j + \epsilon \tag{6.3.7}$$

$$\leq 2\epsilon. \tag{6.3.8}$$

The argument for  $p_l$  is exactly the same.  $\square$

## 6.4 Testable learning via moment matching

### 6.4.1 Warm-up: testable learning over the hypercube via $k$ -wise independence

The main ideas of our approach are already illustrated in the setting of the Boolean hypercube. A key technical ingredient for us will be the fact that an almost  $k$ -wise independent distribution is statistically close to being truly  $k$ -wise independent (Theorem 6.2.6).

**Theorem 6.4.1.** *Let  $\mathcal{C}$  be any concept class that is  $\frac{\epsilon}{4}$ -fooled by  $k$ -wise independence. Then  $\mathcal{C}$  can be testably learned w.r.t.  $\text{Unif}^d$  up to excess error  $\epsilon$  with time and sample complexity  $d^{O(k)}/\epsilon^2$ .*

*Proof.* Let the unknown labeled distribution be  $D$ . Let  $S \subseteq D^m$  be the labeled sample given to  $(T, A)$  (where the sample size  $m$  will be picked later), and let  $S_X$  be its (unlabeled) marginal. Let  $\hat{D}_m$  be the induced empirical distribution, i.e. the uniform distribution over  $S_X$ .

The tester  $T$  and algorithm  $A$  are simple: the tester checks that the empirical moments (or biases) up to degree  $k$  are all no larger than  $\delta = \epsilon d^{-k}/4$  in magnitude (i.e. that the empirical distribution is  $(\delta, k)$ -independent), and the algorithm runs degree- $k$  polynomial regression (Theorem 6.2.2) over the sample.

It is clear that when  $D$  indeed has marginal exactly  $\text{Unif}^d$  (or indeed any  $(\delta/2, k)$ -independent distribution), then by taking  $m = d^k/\delta^2 = d^{k+2}/\epsilon^2$  sufficiently large, we can ensure with high probability all the empirical moments of order at most  $k$  concentrate about their true moments up to  $\delta$  (by a standard Hoeffding plus union bound). That is, with high probability  $\hat{D}_m$  will indeed be  $(\delta, k)$ -independent, and the tester will accept. This verifies completeness.

To verify soundness, suppose that  $\hat{D}_m$  is indeed  $(\delta, k)$ -independent. By Theorem 6.2.6, this means that  $\hat{D}_m$  has TV distance at most  $\delta d^k = \epsilon/4$  from a truly

$k$ -wise independent distribution. This in turn means that  $\widehat{D}_m$  (and indeed any  $(\delta, k)$ -independent distribution)  $\frac{\epsilon}{2}$ -fools  $\mathcal{C}$ . We now appeal to duality, stated here in some generality as Theorem 6.3.2, although in the setting of the hypercube this theorem reduces exactly to the form in Theorem 6.2.5. Formally, observe that for every  $f \in \mathcal{C}$ , condition (a) of Theorem 6.3.2 is satisfied (with  $D = \text{Unif } f \in \mathcal{G}^d$ , and where the slack parameter  $\Delta$  is now simply  $\delta$  in every coordinate). This allows us to apply Corollary 6.3.3 to conclude that there exist  $\epsilon$ -sandwiching polynomials for  $\mathcal{C}$  w.r.t.  $\widehat{D}_m$ . By Theorem 6.2.2, this ensures the learner succeeds at learning  $\mathcal{C}$  up to error  $\text{opt}(\mathcal{C}, D) + \epsilon$  with high probability. This proves the theorem.  $\square$

We may apply this theorem to obtain testable learning w.r.t.  $\text{Unif } f \in \mathcal{G}^d$  for halfspaces, degree-2 PTFs, and constant-depth circuits.

**Corollary 6.4.2.** *Let  $\mathcal{C}$  be the class of halfspaces over  $f \in \mathcal{G}^d$ . Let  $\epsilon > 0$ , and let  $k = \widetilde{O}(1/\epsilon^2)$ . Then  $\mathcal{C}$  is  $\epsilon$ -fooled by  $k$ -wise independence [DGJ<sup>+</sup>10], and hence it can be testably learned w.r.t.  $\text{Unif } f \in \mathcal{G}^d$  up to excess error  $\epsilon$  with time and sample complexity  $d^{O(k)}$ .*

**Corollary 6.4.3.** *Let  $\mathcal{C}$  be the class of degree-2 polynomial threshold functions over  $f \in \mathcal{G}^d$ . Let  $\epsilon > 0$ , and let  $k = \widetilde{O}(1/\epsilon^9)$ . Then  $\mathcal{C}$  is  $\epsilon$ -fooled by  $k$ -wise independence [DKN10], and hence it can be testably learned w.r.t.  $\text{Unif } f \in \mathcal{G}^d$  up to excess error  $\epsilon$  with time and sample complexity  $d^{O(k)}$ .*

**Corollary 6.4.4.** *Let  $\mathcal{C}$  be the class of depth- $t$   $\text{AC}^0$  circuits of size  $s$  over  $f \in \mathcal{G}^d$ . Let  $\epsilon > 0$ , and let  $k = (\log s)^{O(t)} \log(1/\epsilon)$ . Then  $\mathcal{C}$  is  $\epsilon$ -fooled by  $k$ -wise independence [Bra10, Tal17, HS19], and hence it can be testably learned w.r.t.  $\text{Unif } f \in \mathcal{G}^d$  up to excess error  $\epsilon$  with time and sample complexity  $d^{O(k)}$ .*

## 6.4.2 A general algorithm using moment matching

We now give a more general algorithm for testable learning that does not need the target distribution to be  $k$ -wise independent. In this case, our tester will

check that the low-degree moments of the empirical distribution are close to those of the target distribution. The correctness of our tester is a consequence of duality (Theorem 6.3.2).

**Theorem 6.4.5.** *Let  $D_X$  be a distribution on  $X$ , and let  $\mathcal{C}$  be a concept class mapping  $X$  to  $\mathbb{R}$ . Let  $k \geq 1, \Delta \geq \mathbb{R}_+^{j^{(k,d)}}$  be the degree and slack parameters, as in Definition 6.3.1, and let  $\epsilon > 0$  be the error parameter. Suppose the following conditions hold:*

- (a) *(Empirical moments concentrate around true moments.) There exists  $m$  large enough that with high probability over a sample  $S_X \sim D_X^m$ , the corresponding empirical distribution  $\hat{D}_m$  matches moments of degree at most  $k$  with  $D_X$  up to slack  $\Delta$ .*
- (b) *(Existence of sandwiching polynomials with bounded coefficients for  $\mathcal{C}$ , or equivalently approximate moment matching fools  $\mathcal{C}$ .) For every  $f \in \mathcal{C}$ , there exist degree- $k$  sandwiching polynomials  $p_l \leq f \leq p_u$  such that*

$$\mathbb{E}_{D_X} [p_u - f] + h\Delta, |p_u| \leq \frac{\epsilon}{2}, \quad \mathbb{E}_{D_X} [f - p_l] + h\Delta, |p_l| \leq \frac{\epsilon}{2}.$$

*(Recall that for a degree- $k$  polynomial  $p(x) = \sum_I p_I x_I$ , we use  $h\Delta, |p_I|$  to denote  $\sum_{|I| \leq k} |p_I| \Delta^{|I|}$ .)*

*Equivalently, for every  $f \in \mathcal{C}$  and for any distribution  $D^\theta$  whose moments up to order  $k$  match those of  $D_X$  up to  $\Delta$ , we have  $|\mathbb{E}_{D^\theta}[f] - \mathbb{E}_{D_X}[f]| \leq \frac{\epsilon}{2}$ .*

*Then  $\mathcal{C}$  can be testably learned w.r.t.  $D_X$  up to excess error  $\epsilon$  using time and sample complexity  $m + d^{O(k)}$ . Moreover, the tester  $T$  and learner  $A$  are simple:  $T$  tests whether the empirical moments up to order  $k$  match those of  $D_X$  up to  $\Delta$ , and  $A$  performs degree- $k$  polynomial regression over the sample (Theorem 6.2.2).*

*Proof.* Let the unknown labeled distribution be  $D$ , and let  $S \sim D^m$  be the sample given to  $(T, A)$ . First we verify completeness. By assumption, when  $D$  indeed has



marginal  $D_X$ , then  $m$  is large enough that with high probability over  $S$ , the empirical moments concentrate about the true moments up to  $\Delta$ , and hence  $T$  accepts.

As for soundness, suppose that  $T$  accepts, i.e. that the empirical distribution  $\widehat{D}_m$  indeed matches order- $k$  moments with  $D_X$  up to  $\Delta$ . Observe that our condition (b) is the same as condition (b) of Theorem 6.3.2 is satisfied. Thus we may apply Corollary 6.3.3 (with  $D = D_X$  and  $D^\theta = \widehat{D}_m$ ) to conclude that there exist degree- $k$   $\epsilon$ -sandwiching polynomials for  $\mathcal{C}$  w.r.t.  $\widehat{D}_m$ . By Theorem 6.2.2, we have that degree- $k$  polynomial regression achieves error  $\text{opt}(\mathcal{C}, D) + \epsilon$  with high probability. (This implicitly assumes that the degree- $k$  polynomial fitting  $S$  will generalize to  $D$ , which will be true by classic VC theory whenever  $m \geq d^{O(k)}/\epsilon^2$  since the VC dimension of degree- $k$  polynomials (with bounded coefficients, as here) is at most  $d^{O(k)}$ . If this is not the case, we may replace  $m$  with  $m + d^{O(k)}/\epsilon^2$ .)  $\square$

To see how Theorem 6.1.1 may be recovered from this, simply set  $\Delta$  to be  $\delta$  in every coordinate, and also rescale  $\epsilon$  appropriately.

## 6.5 Testably learning functions of halfspaces over strictly subexponential distributions

In this section we apply Theorem 6.4.5 to prove that we can testably learn functions of halfspaces w.r.t. a target marginal  $D_X$  on  $X = \mathbb{R}^d$  that is anticoncentrated and has strictly subexponential tails in the sense of Definition 6.1.3 from the introduction.

**Definition 6.5.1** (Restatement of Definition 6.1.3). We say a distribution  $D_X$  on  $\mathbb{R}^d$  is anticoncentrated and has  $\alpha$ -strictly subexponential tails if the following hold:

- (a) For all  $k$   $k_2 = 1$ ,  $\mathbb{P}[|j^k x, u| > t] \leq \exp(-C_1 t^{1+\alpha})$  for some constant  $C_1$ .
- (b) For all  $k$   $k_2 = 1$  and  $k \geq 2$ ,  $\mathbb{E}[|j^k x, u|^{1/k}] \leq C_2 k^{1/(1+\alpha)}$  for some constant  $C_2$ .

- (c) For all  $ku_2 = 1$  and continuous intervals  $T \subseteq \mathbb{R}$ , we have  $\mathbb{P}[h(x, u) \in T] \leq C_3/T$  for some constant  $C_3$ .

The first two conditions are a strengthening of the usual definition of subexponential distributions (see e.g. [Ver18]), and standard arguments show that the two are actually equivalent. The third asks directional marginals of  $D_X$  to be anticoncentrated. Examples include all *strongly* logconcave distributions, which satisfy this definition with  $\alpha = 1$  (see e.g. [SW14, §5.1] or [Led01, Thm 2.15]). This class includes the standard Gaussian distribution, the uniform distribution on the unit  $d$ -dimensional sphere and more generally, uniform distribution on any convex body with smooth boundaries (e.g., Gaussian smoothing of arbitrary convex bodies).

Throughout this section, let  $\mathcal{C}$  be the class of functions of  $p$  halfspaces over  $\mathbb{R}^d$ , i.e. functions  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  of the form

$$f(x) = g(\text{sign}(hw^1, x) + \theta_1, \dots, \text{sign}(hw^p, x) + \theta_p) \quad (6.5.1)$$

for some  $w^1, \dots, w^p \in \mathbb{R}^d$  (where we use superscripts to avoid confusion with coordinate notation),  $\theta_1, \dots, \theta_p \in \mathbb{R}$ , and  $g : \{-1, 1\}^p \rightarrow \mathbb{R}$ . We focus on the setting where  $p$  is a constant. Also let  $D_X$  be some fixed distribution that is anticoncentrated and  $\alpha$ -strictly subexponential. We will prove the following theorem, stated earlier as Theorem 6.1.4.

**Theorem 6.5.2.** *Let  $\mathcal{C}$  be the class of functions of  $p$  halfspaces over  $\mathbb{R}^d$ , as above. Assume that  $p = O(1)$ . Let  $D_X$  be a distribution that is anticoncentrated and  $\alpha$ -strictly subexponential. Then  $\mathcal{C}$  can be testably learned w.r.t.  $D_X$  up to excess error  $\epsilon$  using  $d^{\tilde{O}(\epsilon^{-(1+\alpha)/\alpha})}$  sample and time complexity.*

In particular whenever  $\alpha = 1$ , as for strongly logconcave distributions including  $N(0, I_d)$ , we obtain a  $d^{\tilde{O}(1/\epsilon^2)}$ -time algorithm.

We now describe our proof plan. To use Theorem 6.4.5, we must show that approximately matching the low-degree moments of  $D_X$  fools functions of halfspaces.

Work due to [KM13] introduced an argument for this problem based on general techniques from the classical theory of moments and the method of metric distances in probability [KR96, RKSF13]. Their broad proof approach was to use [KR96, Thm 2] to show that closeness in moments of two distributions implies closeness in the  $\lambda$ -distance (Definition 6.5.3), and then relate this to the CDF distance, which directly relates to fooling halfspaces. For our purposes, a direct application of [KR96, Thm 2] does not suffice. Instead, we directly analyze the  $\lambda$ -distance under the assumption that the moments of  $D_X$  grow in a strictly subexponential fashion. We begin with the technical lemmas we need, and then prove Theorem 6.5.2 in the final subsection.

### 6.5.1 Moment closeness implies distribution closeness

**Definition 6.5.3** ( $\lambda$ -distance, see e.g. [Zol84], [RKSF13, Chap 10]). For a distribution  $P$  on  $\mathbb{R}^p$ , let  $\phi_P : \mathbb{R}^p \rightarrow \mathbb{C}$  given by  $\phi_P(t) = \mathbb{E}_z \sim P[e^{i\langle t, x \rangle}]$  be its characteristic function. For two distributions  $P, P^\theta$  on  $\mathbb{R}^p$ , define the  $\lambda$ -distance between them as follows:

$$d_\lambda(P, P^\theta) = \min_{T>0} \max \left\{ \max_{|t| \leq T} |\phi_P(t) - \phi_{P^\theta}(t)|, \frac{1}{T} \right\}.$$

We prove an approximate version of [KR96, Thm 1] (see also [RKSF13, Thm 10.3.4]), bounding the  $\lambda$ -distance between two distributions whose moments approximately match and grow with the degree  $k$  in a strictly subexponential fashion, i.e. as  $k^{k/(1+\alpha)}$ .

**Lemma 6.5.4.** *Let  $k \geq 2 \in \mathbb{N}$  be even. Let  $P$  be a distribution on  $\mathbb{R}^p$  such that for all  $k \leq k \leq 1$ ,*

$$\mathbb{E}_z \sim P[j^k h u, z^j] \leq M_k := p^{k/2} C_2^k k^{k/(1+\alpha)}$$

*for some constant  $C_2$ . Let  $P^\theta$  be a distribution that approximately matches moments up to order  $k$  with  $P$  in the following strong sense: for all  $j \leq k$  and  $k \leq k \leq 1$ ,*

$$\left| \mathbb{E}_z \sim P[h u, z^j] - \mathbb{E}_{z^\theta} \sim P^\theta[h u, z^j] \right| \leq \eta_j := \frac{j!}{2^k} \left( \frac{6M_k}{k!} \right)^{(j+1)/(k+1)} = \frac{j!}{2^k} \left( \frac{p}{C_4 k^\alpha} \right)^{j+1}$$

for some constant  $C_4$  depending only on  $C_2$ . Then

$$d_\lambda(P, P^\theta) \leq C_4 \bar{\rho} k^{-\alpha/(1+\alpha)}.$$

*Proof.* To control  $d_\lambda(P, P^\theta)$ , we need to control  $\max_{t \in T} |f(\phi_P(t)) - \phi_{P^\theta}(t)|$  as a function of  $T$ . To this end, fix any direction  $u \in \mathbb{R}^p$  with  $\|u\| = 1$ , and let  $t = \tau u$  for  $\tau \in [0, T]$  be a vector in that direction satisfying  $\|t\| \leq T$ . Let  $\phi_1(\tau) = \phi_P(\tau u)$  and  $\phi_2(\tau) = \phi_{P^\theta}(\tau u)$  be the characteristic functions of  $P$  and  $P^\theta$  along  $u$ . We may Taylor expand  $\phi_1 - \phi_2$  up to degree  $k$  as follows:

$$\phi_1(\tau) - \phi_2(\tau) = \sum_{0 \leq j < k} \frac{\phi_1^{(j)}(0) - \phi_2^{(j)}(0)}{j!} \tau^j + \frac{\phi_1^{(k)}(\tau^\theta) - \phi_2^{(k)}(\tau^\theta)}{k!} \tau^k \quad (6.5.2)$$

for some  $\tau^\theta \in [0, \tau]$ .

The crucial fact we use now is that the derivatives of the characteristic function encode its moments. Indeed, for any  $\tau$ ,

$$\phi_1(\tau) = \mathbb{E}_D[e^{i\tau hz, u}] \Rightarrow \phi_1^{(j)}(\tau) = \mathbb{E}_P[i^j hz, u]^j e^{i\tau hz, u},$$

so that in particular  $j\phi_1^{(j)}(0) = j\mathbb{E}_P[hz, u]^j$  for all  $j$  (and similarly for  $\phi_2$ ). This means  $j\phi_1^{(0)}(0) - \phi_2^{(0)}(0) = 0$ , and for each  $1 \leq j < k$ , by our assumption that  $P^\theta$  approximately moment matches  $P$ , we have

$$j\phi_1^{(j)}(0) - \phi_2^{(j)}(0) \leq \eta_j. \quad (6.5.3)$$

At degree  $k$ , we have

$$j\phi_1^{(k)}(\tau^\theta) = j\mathbb{E}_P[i^k hz, u]^k e^{i\tau^\theta hz, u} = \mathbb{E}_P[j^k hz, u]^k = M_k. \quad (6.5.4)$$

And since  $\mathbb{E}_{P^\theta}[hz^\theta, u]^k = \mathbb{E}_P[hz, u]^k + \eta_k$ , we similarly have

$$j\phi_2^{(k)}(\tau^\theta) = M_k + \eta_k \leq 2M_k \quad (6.5.5)$$

Substituting Eqs. (6.5.3) to (6.5.5) into Eq. (6.5.2), we obtain

$$|\phi_1(\tau) - \phi_2(\tau)| \leq \sum_{1 \leq j < k} \frac{\eta_j}{j!} \tau^j + \frac{3M_k}{k!} \tau^k =: F(\tau), \quad (6.5.6)$$

where we have denoted the expression on the RHS by  $F(\tau)$  for convenience. Since  $F(\tau)$  is clearly increasing in  $\tau$  and independent of  $u$ , we have  $\max_{t \leq T} \int \phi_P(t) \phi_{P^0}(t) j g < F(T)$ . This means that

$$d_\lambda(P, P^0) \leq \max_{t \leq T} \int \phi_P(t) \phi_{P^0}(t) j g, \frac{1}{T} g \leq \max_{t \leq T} F(T), \frac{1}{T} g,$$

and our job now is to pick  $T > 0$  that minimizes the RHS.

This is equivalent to picking the largest  $T$  such that  $F(T) \leq \frac{1}{T}$ , i.e.

$$TF(T) = \sum_{j=1}^k \frac{\eta_j}{j!} T^{j+1} + \frac{3M_k}{k!} T^{k+1} \leq 1.$$

Let us divide this further into two sufficient conditions:

$$\sum_{j=1}^k \frac{\eta_j}{j!} T^{j+1} \leq \frac{1}{2} \quad \text{and} \quad \frac{3M_k}{k!} T^{k+1} \leq \frac{1}{2}.$$

The second condition is equivalent to

$$T = \left( \frac{k!}{6M_k} \right)^{1/(k+1)} = \left( \frac{k!}{6p^{k/2} C_2^k k^{k/(1+\alpha)}} \right)^{1/(k+1)} \asymp \frac{k^{\alpha/(1+\alpha)}}{p^{\frac{1}{p}}},$$

by Stirling's approximation. As for the first, we have picked  $\eta_j$  exactly such that when we plug in this value of  $T$ , for each  $1 \leq j < k$  we have

$$\eta_j = \frac{j!}{2k} T^{-(j+1)} = \frac{j!}{2k} \left( \frac{6M_k}{k!} \right)^{(j+1)/(k+1)} \Rightarrow \frac{\eta_j}{j!} T^{j+1} = \frac{1}{2k}.$$

Summing over  $1 \leq j < k$  verifies the first condition. Thus for this  $T$ , we have

$$d_\lambda(P, P^0) \leq \max_{t \leq T} F(T), \frac{1}{T} g \leq \frac{1}{T} \cdot \frac{1}{p^{\frac{1}{p}} k^{-\alpha/(1+\alpha)}},$$

proving the lemma.  $\square$

We offer some remarks to guide the reader through these calculations. For our application,  $P$  will be the distribution of  $\|x, v\|$  for  $x \in D_X$  and  $\|v\| = \rho_{\bar{p}}$ . The key idea is simply to use a Taylor approximation of the  $\lambda$ -distance to reduce the issue to one of moment closeness. The final calculation amounts to solving for  $T$  satisfying

$T^{k+1} \cap \frac{k!}{M_k}$ , where the denominator is the  $k^{\text{th}}$  moment of  $P$ . We then set each  $\eta_j$  small enough to make the lower degree terms minor. Note that the  $j^{\text{th}}$  moment of  $P$  scales as  $M_j = p^{j/2} C_2^j j^{j/(1+\alpha)}$ , and we have  $M_k^{j/k} = p^{j/2} C_2^j k^{j/(1+\alpha)}$ . Thus loosely speaking, the slack  $\eta_j$  may be viewed in relative terms as follows:

$$\frac{\eta_j}{M_j} = \frac{1}{2k} \frac{j!}{M_j} \left( \frac{6M_k}{k!} \right)^{(j+1)/(k+1)} \quad (6.5.7)$$

$$\frac{1}{2k} \frac{j!}{M_j} \left( \frac{6M_k}{k!} \right)^{j/k} \quad (6.5.8)$$

$$\frac{1}{2k} \frac{j^j}{p^{j/2} C_2^j j^{j/(1+\alpha)}} \frac{p^{j/2} C_2^j k^{j/(1+\alpha)}}{k^j} \quad (6.5.9)$$

$$\frac{1}{2k} \left( \frac{j}{k} \right)^{j\alpha/(1+\alpha)}. \quad (6.5.10)$$

This relative slack factor is only about  $1/\text{poly}(k)$  for  $j = O(1)$  but for  $j = \Theta(k)$  it becomes  $\exp(-\Theta(k))$ , which seems unavoidable with our method. Finally, note also that if  $P$ 's moments scaled only as a subexponential instead of a strictly subexponential distribution, i.e. if the  $k^{\text{th}}$  moment of  $P$  scaled with  $k$  as  $k^k$ , then the key calculation for  $T$  becomes vacuous. More involved techniques (see e.g. [RKSF13, Thm 10.3.1]) still have something to say in this situation when certain stricter moment conditions hold, but the direct Taylor expansion approach fails.

The following lemma is a convenient distillation of the rest of the argument from [KM13], where the  $\lambda$ -distance is related to the Levy distance (using [Gab81]), which in turn is related to the CDF distance (using anticoncentration), and which leads finally to the desired conclusion.

**Lemma 6.5.5** (Implicit in [KM13, §3.3]). *Let  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  be a function of  $p$  halfspaces as above, and also let distributions  $D, D^\theta$  on  $\mathbb{R}^d$  and  $P, P^\theta$  on  $\mathbb{R}^p$  be as above. Assume that for any continuous interval  $T \subseteq \mathbb{R}$ , each coordinate  $z_j$  of  $z \sim P$  satisfies  $\mathbb{P}[z_j \in T] \geq \Theta(jT)$ . Suppose that  $d_\lambda(P, P^\theta) \leq \delta$ . Let  $N(\delta)$  be such that  $\mathbb{P}_{z \sim P}[kz^\theta k_1 > N(\delta)] \leq \delta$  and  $\mathbb{P}_{z^\theta \sim P^\theta}[kz^\theta k_1 > N(\delta)] \leq \delta$ . Then*

$$j \mathbb{E}_D[f] - \mathbb{E}_{D^\theta}[f] \leq O(2^p \delta (\log N(\delta) + 2 \log(1/\delta))^p).$$

### 6.5.2 Approximate low-degree moment matching fools functions of half-spaces

We now prove our main structural result, which is that any distribution that approximately matches the low-degree moments of  $D_X$  fools functions of halfspaces.

**Theorem 6.5.6.** *Let  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  be of the form in Eq. (6.5.1). Let  $D_X$  be a distribution that is anticoncentrated and  $\alpha$ -strictly subexponential. For any  $k \geq 2\mathbb{N}$ , let  $\Delta \geq \mathbb{R}_+^{j \in [k, d]}$  be such that for each  $I \in [k, d]$  with  $|I| = j$ ,*

$$\Delta_I = \frac{\rho_{\bar{p}} j!}{2k} \left( \frac{1}{C_4 k^{\alpha/(1+\alpha)}} \right)^{j+1} \quad (6.5.11)$$

for some constant  $C_4 > 0$ . Then for any distribution  $D^\theta$  whose moments up to order  $k$  match those of  $D_X$  up to  $\Delta$ , we have

$$j \mathbb{E}_D[f] - \mathbb{E}_{D^\theta}[f] \leq k^{-\alpha/(1+\alpha)} \rho_{\bar{p}} (C \log(\rho_{\bar{p}} k^{\alpha/(1+\alpha)}))^{2p}$$

for some constant  $C$ . In particular, for  $p = O(1)$ , we have  $j \mathbb{E}_D[f] - \mathbb{E}_{D^\theta}[f] \leq \tilde{O}(k^{-\alpha/(1+\alpha)})$ .

*Proof.* Let  $D = D_X$ . Assume without loss of generality that  $w^1, \dots, w^p$  are unit vectors, and let  $W \in \mathbb{R}^{p \times d}$  be the matrix with the  $w^i$  as its rows. Let  $P$  be the distribution (on  $\mathbb{R}^p$ ) of  $Wx$  for  $x \sim D$ , and define  $P^\theta$  similarly. We would like to apply Lemma 6.5.4 to  $P$  and  $P^\theta$ . To do so, we must first verify moment closeness. Let  $u \in \mathbb{R}^p$  be a unit vector, and let  $v = W^T u \in \mathbb{R}^d$ . For any multi-index  $I \in \mathbb{N}^d$ , let  $v_I$  denote  $\prod_{j \in [d]} v_j^{I_j}$ . Then for any  $j$ ,

$$\mathbb{E}_P[hz, u^j] = \mathbb{E}_D[hx, v^j] \quad (6.5.12)$$

$$= \mathbb{E}_D \left[ \sum_{|I|=j} x_I v_I \right] \quad (6.5.13)$$

$$= \sum_{|I|=j} v_I \mathbb{E}_D[x_I]. \quad (6.5.14)$$

We place a crude upper bound on each  $jv_Ij$  as follows. Since  $W$  has Frobenius norm  $kWk_F = \rho_{\bar{p}}$ , we have

$$kvk_{\gamma} \quad kvk \quad kukkWk_F \quad \rho_{\bar{p}},$$

so that in particular for each  $I$  with  $|I| = j$ ,  $jv_Ij = \prod_{i \in I} jv_i^{I_i} \quad kvk_{\gamma}^{jI} \quad p^{j/2}$ . Thus

$$j \mathbb{E}_P[hz, u^{I^j}] \quad z^0 \mathbb{E}_{P^0}[hz^0, u^{I^j}] = j \mathbb{E}_D[hx, v^{I^j}] \quad x^0 \mathbb{E}_{D^0}[hx^0, v^{I^j}] \quad (6.5.15)$$

$$\sum_{|I|=j} jv_Ij \mathbb{E}_D[x_I] \quad \mathbb{E}_{D^0}[x_I] \quad (6.5.16)$$

$$d^j p^{j/2} \sup_{|I|=j} \Delta_I \quad (6.5.17)$$

$$\eta_j, \quad (6.5.18)$$

where  $\eta_j$  is as defined in Lemma 6.5.4, and the final inequality follows since we have picked  $\Delta$  in the theorem statement precisely such that  $\sup_{|I|=j} \Delta_I = d^j p^{j/2} \eta_j$ . Also observe that

$$\mathbb{E}_P[hu, z^{I^k}] = \mathbb{E}_D[hx, v^{I^k}] = kvk^k C_2^k k^{k/(1+\alpha)} \quad p^{k/2} C_2^k k^{k/(1+\alpha)}$$

Now we apply Lemma 6.5.4 to conclude that  $d_{\lambda}(P, P^0) \leq \rho_{\bar{p}} k^{-\alpha/(1+\alpha)}$ .

To finish the proof, we appeal to Lemma 6.5.5. For this we must first verify anticoncentration of  $P$  and estimate  $N(\delta)$  as defined in that lemma. Observe first that for any  $i \in [d]$ , the  $i^{\text{th}}$  coordinate of  $z \sim P$  (resp.  $z^0 \sim P^0$ ) is precisely  $hw^i, x^i$  for  $x \sim D$  (resp.  $hw^i, x^0$  for  $x^0 \sim D^0$ ). Anticoncentration of each coordinate of  $z$  follows immediately from Definition 6.5.1(c). To estimate  $N(\delta)$ , we will use a simple Chebyshev-style bound. For any coordinate  $i \in [d]$  and any even degree  $j \leq k$ , we have

$$\mathbb{P}_D[jhw^i, x^i] > t \leq \frac{\mathbb{E}_D[hw^i, x^i]^j}{t^j} \leq \frac{C_2^j j^{j/(1+\alpha)}}{t^j}.$$

And since  $D^0$  approximately matches moments with  $D$ , by a similar calculation as earlier (now with  $kw^i k_{\gamma} \leq 1$  in place of  $kvk_{\gamma} \leq \rho_{\bar{p}}$ ),

$$\mathbb{E}_{D^0}[hw^i, x^0] \leq \mathbb{E}_D[hw^i, x^i] + \sum_{|I|=j} jv_Ij \Delta_I \leq \mathbb{E}_D[hw^i, x^i] + \eta_j/p^{j/2} \leq 2 \mathbb{E}_D[hw^i, x^i],$$



and so

$$\mathbb{P}_{D^\theta}[\sum_{i=1}^j hw^i, x^\theta / j > t] = \frac{\mathbb{E}_{D^\theta}[\sum_{i=1}^j hw^i, x^\theta / j]}{t^j} = \frac{2C_2^j j^{j/(1+\alpha)}}{t^j} = 2 \left( \frac{C_2 j^{1/(1+\alpha)}}{t} \right)^j.$$

We need  $t$  such that the RHS is at most  $\delta/p$ . For this it suffices to set  $j = 2 \log(p/\delta)$  and  $t = C_2 j^{1/(1+\alpha)}$  for this. By a union bound over the  $p$  coordinates of  $z \in P$  (similarly  $z^\theta \in P^\theta$ ), we see that we may take  $N(\delta) = t = O((\log(p/\delta))^{1/(1+\alpha)})$ .

We are now ready to apply Lemma 6.5.5 with this  $N(\delta)$  and  $\delta \wedge \sqrt[p]{k}^{-\alpha/(1+\alpha)}$ . Substituting these expressions in, we get that

$$\begin{aligned} \mathbb{E}_D[f] - \mathbb{E}_{D^\theta}[f] & \leq O(2^p \delta (\log N(\delta) + 2 \log(1/\delta))^p) & (6.5.19) \\ & \leq 2^p \delta \left( C^\theta \log\left(\frac{1}{\delta} \log \frac{p}{\delta}\right) \right)^p & \text{(for some constant } C^\theta > 0) \\ & \leq \delta \left( C \log\left(\frac{p}{\delta}\right) \right)^{2p} & \text{(for some constant } C > 0) \\ & \leq k^{-\alpha/(1+\alpha)} \sqrt[p]{k}^{-\alpha/(1+\alpha)} \left( C \log\left(\sqrt[p]{k}^{-\alpha/(1+\alpha)}\right) \right)^{2p}, & (6.5.20) \end{aligned}$$

as claimed.  $\square$

We pause to note an interesting corollary of this theorem, stated informally earlier as Corollary 6.1.5, which states that any  $D^\theta$  that approximately matches low-degree moments with  $D_X$  must be anticoncentrated.

**Corollary 6.5.7.** *Let  $\epsilon > 0$ ,  $k = \tilde{O}(\epsilon^{-(1+\alpha)/\alpha})$ , and  $\Delta$  be as in Theorem 6.5.6, with  $p = 2$ . Let  $D^\theta$  be any distribution whose moments up to order  $k$  match those of  $D_X$  up to  $\Delta$ . Then for any  $k$ -bit  $u$  and any continuous interval  $T \subseteq \mathbb{R}$ ,  $|\mathbb{P}_{D^\theta}[\sum_{i=1}^k u_i \in T] - \mathbb{P}_{D_X}[\sum_{i=1}^k u_i \in T]| \leq \epsilon + \Theta(kT)$ .*

*Proof.* Write  $T = [\theta, \theta']$  for some  $\theta < \theta' \in \mathbb{R}$ , and consider the function  $f(x) = \text{sign}(\sum_{i=1}^k u_i - \theta) \wedge \text{sign}(\theta' - \sum_{i=1}^k u_i)$  (where  $b_1 \wedge b_2 = 1$  iff  $b_1 = b_2 = 1$ ). Clearly  $\sum_{i=1}^k u_i \in T$  iff  $f(x) = 1$ . But  $f$  is an intersection of two halfspaces, and we know by Theorem 6.5.6 that  $|\mathbb{E}_{D^\theta}[f] - \mathbb{E}_{D_X}[f]| \leq \epsilon$ . Since  $\mathbb{E}_{D_X}[f] = C_3/T$  by Definition 6.5.1(c), the statement follows.  $\square$

In fact, the same reasoning tells us that for any collection of  $p = O(1)$  intervals  $T$  and any  $\|k\| = 1$ ,  $\mathbb{P}_{D^0}[\|k\| \geq T] = \mathbb{P}_{D_X}[\|k\| \geq T] + \epsilon$ .

### 6.5.3 Proof of Theorem 6.5.2

The final ingredient for the proof of Theorem 6.5.2 is the following lemma, which gives a bound on the sample complexity required for the empirical moments of  $D_X$  to concentrate about their true moments.

**Lemma 6.5.8.** *Let the degree parameter be  $k$ , and the slack parameter  $\Delta \geq \mathbb{R}_+^{j \mid (k,d)j}$  be as in Theorem 6.5.6. Assume  $p = O(1)$ . Then drawing a sample of size  $m = d^{\tilde{O}(k)}$  from  $D_X$  is sufficient to ensure that with high probability, the empirical moments of order at most  $k$  match those of  $D_X$  up to slack  $\Delta$ .*

*Proof.* Let  $D = D_X$ , and let  $\hat{D}_m$  denote the empirical distribution on a sample  $S$  of size  $m$  drawn from  $D$ . We would like to ensure that for every  $I \in \mathcal{I}(k, d)$ ,  $j \in \hat{D}_m[x_I] - \mathbb{E}_D[x_I] \leq \Delta_I$ . It suffices to consider the case when  $x_I$  has the weakest concentration, and this is clearly when  $|I| = k$  and in fact when  $I = (k, 0, \dots, 0)$  (without loss of generality), so that  $x_I = x_1^k$ .

Let  $Z$  denote the random variable  $\mathbb{E}_{\hat{D}_m}[x_1^k] - \mathbb{E}_D[x_1^k]$ . For our purposes it is sufficient to use a crude Chebyshev bound, although higher moment analogs will give a slightly better bound. We have  $\text{Var}[Z] = \frac{1}{m} \text{Var}[x_1^k] = \frac{1}{m} C_2^{2k} (2k)^{2k/(1+\alpha)}$ . Thus

$$\mathbb{P}[Z > \Delta_I] \leq \frac{\text{Var}[Z]}{\Delta_I^2} \tag{6.5.21}$$

$$\leq \frac{1}{m} \frac{C_2^{2k} (2k)^{2k/(1+\alpha)}}{\Delta_I^2} \tag{6.5.22}$$

$$\leq \frac{1}{m} C_2^{2k} (2k)^{2k/(1+\alpha)} \left( \frac{2k}{p} \frac{d^k}{k!} \right)^2 (C_4 k^{\alpha/(1+\alpha)})^{2(k+1)} \tag{6.5.23}$$

$$\leq \frac{1}{m} (dk)^{O(k)}, \tag{6.5.24}$$

after plugging in the value of  $\Delta_I$  when  $|I| = k$  from Eq. (6.5.11) and some manipulation. This is at most  $\delta$  if  $m \geq (dk)^{O(k)}/\delta$ .

For moment closeness to hold simultaneously for all  $I \geq I(k, d)$  with high probability, we set  $\delta = \Theta(1/j!(k, d)^j) = d^{-\binom{k}{j}}$  and apply a union bound. For this  $\delta$ ,  $m$  may be simplified to  $d^{\tilde{O}(k)}$ , as desired.  $\square$

We are now ready to prove Theorem 6.5.2 using our general algorithm, Theorem 6.4.5.

*Proof of Theorem 6.5.2.* We need to pick  $m$ ,  $k$  and  $\Delta$  suitably as functions of  $\epsilon$  and  $d$ , and verify that the conditions in Theorem 6.4.5 hold. Let  $\Delta$  be as defined in Theorem 6.5.6. By Lemma 6.5.8, it suffices to take  $m = d^{\tilde{O}(k)}$  to ensure that with high probability, the empirical distribution  $\hat{D}_m$  matches moments of order at most  $k$  with  $\mathcal{N}(0, I_d)$  up to  $\Delta$ . This verifies condition (a) of Theorem 6.4.5. Now for any  $f \in \mathcal{C}$ , we can combine Theorem 6.5.6 with Theorem 6.3.2 to obtain degree- $k$  sandwiching polynomials satisfying condition (b) of Theorem 6.4.5, with  $\epsilon = \tilde{O}(k^{-\alpha/(1+\alpha)})$ , or equivalently  $k = \tilde{O}(\epsilon^{-(1+\alpha)/\alpha})$ . Applying Theorem 6.4.5 completes the proof.  $\square$

## 6.6 Sample complexity of testable learning

In this section we show that the sample complexity of testably learning a class is characterized by its Rademacher complexity. Throughout this section, let  $D_X$  be the target marginal, and let  $\mathcal{C}$  be the concept class (mapping  $X$  to  $\mathbb{R}$ ) that we wish to testably learn w.r.t.  $D_X$ . We remind the reader that we use the term “with high probability” to mean “with probability at least 0.99” for simplicity.

### 6.6.1 Upper bound

We begin with the upper bound, which is essentially just the observation that the empirical Rademacher complexity provides a generalization guarantee that can be estimated to high accuracy from the sample itself. Note that this is an information-theoretic upper bound.

**Theorem 6.6.1.** *Let  $\epsilon > 0$ , and let  $m^\circ$  be such that  $R_{m^\circ}(\mathcal{C}, D_X) \leq \frac{\epsilon}{5}$ . Then  $\mathcal{C}$  can be testably learned w.r.t.  $D_X$  up to excess error  $\epsilon$  with sample complexity  $m^\circ + O(1/\epsilon^2)$ .*

*Proof.* Let  $m = m^\circ + O(1/\epsilon^2)$ . Let  $S = \{(x_i, y_i)\}_{i \in [m]} \in D^m$  be a sample of  $m$  labeled points drawn from  $D$ , and let  $S_X$  denote  $\{x_i\}_{i \in [m]}$ . Our tester  $T$  accepts iff  $\widehat{R}_m(\mathcal{C}, S_X) \leq \frac{\epsilon}{4}$ . Whenever the tester accepts, the learner  $A$  simply performs ERM over  $S$  w.r.t.  $\mathcal{C}$ .

To see why completeness is satisfied, suppose that the true marginal is in fact  $D_X$ . Since  $R_m(\mathcal{C}, D_X) \leq R_{m^\circ}(\mathcal{C}, D_X) \leq \frac{\epsilon}{5}$ , by Eq. (6.2.3) we can ensure that with high probability over  $S_X$ ,  $\widehat{R}_m(\mathcal{C}, S_X) \leq \frac{\epsilon}{4}$ , and so  $T$  will accept.

Soundness holds by a standard argument showing generalization using uniform convergence. Formally, suppose that the tester accepts  $S_X$ , i.e.  $\widehat{R}_m(\mathcal{C}, S_X) \leq \frac{\epsilon}{4}$ . Consider an ERM hypothesis

$$\widehat{f}_m = \arg \min_{f \in \mathcal{C}} \frac{1}{m} \sum_{i \in [m]} \ell(f(x_i), y_i)$$

as well as an optimal hypothesis

$$f^* = \arg \min_{f \in \mathcal{C}} \mathbb{E}_{(x,y) \sim D} [\ell(f(x), y)].$$

Then with high probability over  $S$  we have

$$\begin{aligned} L(\widehat{f}_m) &= \widehat{L}_m(\widehat{f}_m) + \widehat{R}_m(\mathcal{C}) + O\left(\sqrt{\frac{1}{m}}\right) && \text{(by Eq. (6.2.5))} \\ &= \widehat{L}_m(f^*) + \widehat{R}_m(\mathcal{C}) + O\left(\sqrt{\frac{1}{m}}\right) && (\widehat{f}_m \text{ is an ERM hypothesis}) \\ &= L(f^*) + 2\widehat{R}_m(\mathcal{C}) + O\left(\sqrt{\frac{1}{m}}\right) && \text{(by Eq. (6.2.5) again)} \\ &= L(f^*) + \frac{\epsilon}{2} + \frac{\epsilon}{2} = L(f^*) + \epsilon, && (6.6.1) \end{aligned}$$

by our choice of  $m$ . This proves the theorem.  $\square$

### 6.6.2 Lower bound

Now we state the lower bound, which matches the upper bound up to a quadratic factor. Our lower bound can be viewed as a generalization of the argument of [RV22], who proved lower bounds for testable learning for the specific cases of convex sets and monotone functions. We obtain the full range of lower-bounds for any value of  $\epsilon$  purely in terms of Rademacher complexity. The idea here is that no tester with bounded sample complexity  $m$  can distinguish between a distribution  $D_X$  and the uniform distribution on a sufficiently large sample of size  $M - m$  drawn from  $D_X$ . If the Rademacher complexity w.r.t.  $D_X$  at sample size  $M$  is somewhat large, then the large sample can be labeled randomly and still admit nontrivial optimal error, but of course the learner cannot do well on unseen data.

**Theorem 6.6.2.** *Let  $\epsilon > 0$ , and let  $M$  be such that  $R_M(\mathcal{C}, D_X) \leq 5\epsilon$ . Assume that  $M = \Theta(1/\epsilon^2)$  is sufficiently large. Then testably learning  $\mathcal{C}$  up to excess error  $\epsilon$  requires sample complexity at least  $\Omega(\frac{\rho}{M})$ .*

*Proof.* Suppose we had a tester-learner  $(T, A)$  requiring sample complexity only  $m$  where  $m = \frac{\rho}{100}$ . We will show how to “fool”  $(T, A)$  into failing its guarantee by constructing a (random) labeled distribution  $D$  such that (with high probability over the randomness of our construction),

- (a)  $\text{opt}(D, \mathcal{C}) \leq \frac{1}{2} + 2\epsilon$ ;
- (b) with high probability,  $T$  will accept a sample of size  $m$  drawn from  $D$ ; and yet
- (c) with high probability,  $A$ 's output will have error greater than  $\frac{1}{2} - \epsilon$  on  $D$ .

For such a  $D$ , it is clear that the tester-learner pair  $(T, A)$  fails its guarantee in that with high probability, despite  $T$  accepting,  $A$  cannot produce a hypothesis with error at most  $\text{opt}(D, \mathcal{C}) + \epsilon$ .

We construct  $D$  as follows. Draw a sample of  $M$  randomly labeled points  $S = \{(x_i, y_i)\}_{i=1}^M$  ( $D_X = \text{Unif}(\mathcal{X})^M$ ), and let  $S_X$  denote  $\{x_i\}_{i=1}^M$ . Define  $D$  to be the uniform distribution over  $S$ . We now show that with high probability over the draw of  $S$  (including its random labeling), the distribution  $D$  satisfies the required properties.

Denote the size- $m$  sample given to  $(T, A)$  by  $S^\theta \sim D^m$ , and let  $S_X^\theta$  denote its marginal. That is,  $S^\theta$  is an iid subsample of  $S$ . Notice that unless we sample the same element of  $S$  twice,  $S_X^\theta$  is distributed exactly as a sample of  $m$  points drawn directly from  $D_X$ . In other words, the statistical distance between  $S_X^\theta$  formed in this way vs by drawing  $m$  points directly from  $D_X$  is at most the collision probability, which is  $m^2/M \leq 10^{-4}$ . Thus any tester satisfying completeness must accept  $S^\theta$  with high probability. This gives us property (b).

Let us see why property (a) holds with high probability over  $S$ . The idea is that because  $R_M(\mathcal{C}) \leq \Omega(\epsilon)$ , we expect that there exists a classifier in  $\mathcal{C}$  that achieves error at most  $\frac{1}{2} + \Omega(\epsilon)$  on the randomly labeled sample  $S$ . Indeed, observe that the random labels are exactly equivalent to Rademacher random variables. Assuming that  $M$  is sufficiently large and applying Eqs. (6.2.3) and (6.2.4) successively, we obtain that with high probability over the sample  $S$  (together with the realization of the random labels),

$$\left| R_M(\mathcal{C}) - \sup_{f \in \mathcal{C}} \frac{1}{M} \sum_{i=1}^M y_i f(x_i) \right| \leq \epsilon.$$

In particular, since  $R_M(\mathcal{C}) \leq \frac{1}{2} + 5\epsilon$ , there exists  $f \in \mathcal{C}$  such that  $\frac{1}{M} \sum_{i=1}^M y_i f(x_i) \geq \frac{1}{2} - 4\epsilon$ , or equivalently

$$\frac{1}{M} \sum_{i=1}^M \mathbb{1}[f(x_i) \neq y_i] = \frac{1}{M} \sum_{i=1}^M \frac{1 - y_i f(x_i)}{2} = \frac{1}{2} - \frac{1}{2M} \sum_{i=1}^M y_i f(x_i) \leq \frac{1}{2} + 2\epsilon.$$

In other words,  $\text{opt}(D, \mathcal{C}) \leq \frac{1}{2} + 2\epsilon$ .

For property (c), the idea is that the learner, having only seen a minuscule fraction of the randomly labeled  $D$ , cannot possibly output a hypothesis with error

substantially better than  $\frac{1}{2}$  on all of  $D$ . Formally, observe that any classifier  $h$  that  $A$  outputs is stochastically independent of  $S \cap S^\theta$ . This means that in expectation over  $S$  and the randomness of  $A$ ,

$$\mathbb{P}_{(x,y) \sim D}[h(x) \neq y] \leq \frac{m}{M} + \frac{1}{2} \frac{M - m}{M} = \frac{1}{2} \frac{m}{2M} + \frac{1}{2} \frac{\epsilon}{2},$$

since the fact that  $M = \Theta(1/\epsilon^2)$  and  $m = \Omega(\sqrt{M}/100)$  mean that  $\frac{m}{M} = \frac{1}{100\sqrt{M}} < \epsilon$ . It is clear that for sufficiently large  $M = \Theta(1/\epsilon^2)$ , we will actually have  $\mathbb{P}_{(x,y) \sim D}[h(x) \neq y] > \frac{1}{2} - \epsilon$  with high probability over  $S$  and  $A$ . (Such an argument is also formalized as [RV22, Lemma 25].)

□

Note that this theorem becomes stronger if  $\epsilon$  is taken to be a constant. In particular, if we assume  $R_M(\mathcal{C}, D_X) \leq 0.99$ , then the same argument would actually yield a “fooling distribution”  $D$  such that

- (a)  $\text{opt}(D, \mathcal{C}) \leq 0.01$ ;
- (b) with high probability,  $T$  will accept a sample of size  $m$  drawn from  $D$ ;
- (c) with high probability,  $A$ 's output will have error greater than 0.49 on  $D$ .

This would rule out any tester-learner capable of testably learning up to error sufficient to distinguish the case where  $\text{opt}(D, \mathcal{C}) = 0.01$  from  $\text{opt}(D, \mathcal{C}) > 0.49$  (e.g., one with final error guarantee  $10 \cdot \text{opt}(D, \mathcal{C}) + 0.1$ ).

We also give the following stronger version of this lower bound, stated in terms of the behavior of the empirical Rademacher complexity (which is a random variable depending on the sample). This is a very strong lower bound that holds whenever  $\widehat{R}_M(\mathcal{C}, S_X) = 1$  with high probability, because it yields a fooling distribution  $D$  that is in fact perfectly realizable. As we will see in Section 6.6.2.1, this turns out to apply to convex sets and monotone functions. In a sense, this version is not really about

Rademacher complexity but rather the stronger notion of shattering (except with high probability over a sample, like a distribution-specific version of the VC dimension, albeit stronger than VC entropy). Recall that  $\mathcal{C}$  is said to shatter an unlabeled set  $S_X$  if every possible labeling of  $S_X$  can be achieved by some  $f \in \mathcal{C}$ , or equivalently  $\widehat{R}_M(\mathcal{C}, S_X) = 1$ .

**Theorem 6.6.3.** *Let  $M$  be such that with high probability over a size- $M$  sample  $S_X \sim D_X^M$ ,  $\widehat{R}_M(\mathcal{C}, S_X) = 1$ , i.e.  $\mathcal{C}$  shatters  $S_X$ . Consider any tester-learner pair  $(T, A)$  for testably learning  $\mathcal{C}$  up to error sufficient to distinguish the case where  $\text{opt}(D, \mathcal{C}) = 0$  from  $\text{opt}(D, \mathcal{C}) > 0.49$ . Then  $(T, A)$  requires sample complexity at least  $\Omega\left(\frac{\rho}{M}\right)$ .*

*In particular, this rules out any tester-learner with  $\epsilon$ -error guarantee  $\psi(\text{opt}(D, \mathcal{C})) + 0.49$  for any increasing function  $\psi : [0, 1] \rightarrow \mathbb{R}$  satisfying  $\psi(0) = 0$ .*

*Proof.* The proof is a simpler version of the earlier one. Again, suppose we had a tester-learner  $(T, A)$  requiring sample complexity only  $m = \frac{\rho}{100}$ . We construct a labeled distribution  $D$  such that

- (a)  $\text{opt}(D, \mathcal{C}) = 0$ ;
- (b) with high probability,  $T$  will accept a sample of size  $m$  drawn from  $D$ ; and yet
- (c) with high probability,  $A$ 's output will have error greater than 0.49 on  $D$ .

The distribution  $D$  is constructed in exactly the same way: draw a sample of  $M$  randomly labeled points  $S = \{(x_i, y_i) \mid g_i \in \{0, 1\}^M\} \sim (D_X \times \text{Unif}\{0, 1\})^M$ , and define  $D$  to be the uniform distribution over  $S$ . Let  $S_X = \{x_i \mid g_i \in \{0, 1\}^M\}$ . Let  $S^0 \sim D^m$  denote the sample given to  $(T, A)$ , and as before, let us condition on its marginal  $S_X^0$  containing no duplicates (which occurs with high probability).

Property (a) follows immediately from our assumption that with high probability,  $\mathcal{C}$  shatters  $S_X$ . (Note that this also implies that  $S_X$  contains no duplicates.)



Properties (b) and (c) follow by almost exactly the same arguments as before (for the latter, we now use the fact that  $m/M = 1/100$  instead of  $m/M = \epsilon$ ).  $\square$

### 6.6.2.1 Applications

The lower bounds of [RV22] may be viewed as applications of Theorem 6.6.3. The first application is the class of convex sets w.r.t.  $\mathcal{N}(0, I_d)$ , and the second is the class of monotone Boolean functions w.r.t.  $\text{Unif } \mathcal{I}^d$ .

**Theorem 6.6.4** (Implicit in [RV22, Theorem 22]). *Let  $X = \mathbb{R}^d$ ,  $D_X = \mathcal{N}(0, I_d)$ , and  $\mathcal{C}$  be the class of  $\{0, 1\}$ -valued indicator functions of convex sets in  $\mathbb{R}^d$ . Let  $M = 2^{Cd}$  for some small constant  $C > 0$ . Then with probability  $1 - \exp(-\Omega(d))$  over the draw of a size- $M$  sample  $S \sim D_X^M$ ,  $\mathcal{C}$  can shatter  $S$ .*

**Theorem 6.6.5** (Implicit in [RV22, Theorem 23]). *Let  $X = \mathcal{I}^d$ ,  $D_X = \text{Unif } \mathcal{I}^d$ , and  $\mathcal{C}$  be the class of monotone Boolean functions. Let  $M = 2^{Cd}$  for some small constant  $C > 0$ . Then with probability  $1 - \exp(-\Omega(d))$  over the draw of a size- $M$  sample  $S \sim D_X^M$ ,  $\mathcal{C}$  can shatter  $S$ .*

In fact, Rubinfeld and Vasilyan are able to state their lower bounds in a slightly stronger way because of the specific parameters  $M, m$  that these examples above allow. Specifically, for both convex sets and monotone functions, we may take  $M = 2^{\Omega(d)}$ ,  $m = M^{0.01} = 2^{\Omega(d)}$ , and the same argument as in Theorem 6.6.3 can be analyzed more closely to yield a distribution  $D$  such that  $\text{opt}(D, \mathcal{C}) = 0$  and yet the final output of any testable learner with sample complexity  $m$  must have  $\exp(-\Omega(d))$  advantage over random guessing (which is stronger than merely saying the output must have error at least 0.49).

Interestingly, these examples add to what has been called “the emerging analogy between symmetric convex sets in Gaussian space and monotone Boolean functions”; see [DNS22] and references therein.

## 6.7 Discussion

### 6.7.1 Implications for the uniform convergence paradigm

As observed in [RV22], an interesting consequence of the lower bounds in Section 6.6.2.1 is that they demonstrate a strict separation between distribution-specific agnostic learning and testable learning. In the case of both convex sets over  $\mathcal{N}(0, I_d)$  and monotone functions over  $\text{Unif } \mathcal{G}^d$ , Fourier-theoretic arguments are known to give agnostic learners requiring sample complexity only  $2^{\tilde{O}(P_{\bar{d}}/\text{poly}(\epsilon))}$  to learn up to excess error  $\epsilon$  [BT96, KOS08]. In particular, they require only sample complexity  $2^{\tilde{O}(P_{\bar{d}})}$  to learn up to excess error  $\epsilon = 0.1$  (say), which is much smaller than the lower bounds of  $2^{\Omega(d)}$  for testably learning these classes up to  $\epsilon = 0.1$ .

But we have just characterized testable learning in terms of Rademacher complexity, which we know in turn tightly characterizes uniform convergence (Theorem 6.2.8). We draw the following implications from this:

- Uniform convergence is always sufficient for distribution-specific agnostic learning but it is not necessary, as witnessed by the examples of convex sets and monotone functions.
- Uniform convergence is both necessary and sufficient for testable learning.

That is, not only is there a strict separation between distribution-specific agnostic learning and testable learning, it is the latter that is in fact characterized by uniform convergence.

#### **Uniform convergence in distribution-free vs distribution-specific learning.**

In the distribution-free setting, uniform convergence is well-known to be necessary and sufficient for agnostic (as well as realizable) learning, by classic VC theory (see e.g. [SB14, Chapter 6]). Let us clarify that in the distribution-free setting the term “uniform convergence” now means a uniform bound on the generalization gap over

not just all  $f \in \mathcal{C}$  but also all distributions  $D_X$ ; that is, we now care about the *worst-case distribution-free* generalization gap:  $\sup_{D_X} \sup_{f \in \mathcal{C}} |L(f) - \widehat{L}_m(f)|$ . This quantity is tightly governed by the VC-dimension of  $\mathcal{C}$ , a distribution-free, purely combinatorial property.

Meanwhile in the distribution-specific setting, the statistical complexity of agnostic learning is known to be characterized by the  $L^1(D_X)$  metric entropy of the class (aka the log covering number w.r.t. the metric  $\rho(f, g) = \mathbb{E}_{x \sim D_X} [|f(x) - g(x)|]$ , which for  $\{0, 1\}$ -valued  $f, g$  is the same as  $2 \mathbb{P}_{x \sim D_X} [f(x) \neq g(x)]$ ) [BI91]. In this setting, uniform convergence (now in the distribution-specific sense of bounds on  $\sup_{f \in \mathcal{C}} |L(f) - \widehat{L}_m(f)|$  for fixed  $D_X$ ) is sufficient but not necessary. Indeed, we may also view the separations given by convex sets and monotone functions as separations between the metric entropy and the Rademacher complexity of these classes.

A priori, this seems like a surprising difference between distribution-free and distribution-specific agnostic learning. However, one could argue that the more realistic distribution-specific supervised learning model is that of testable learning. Here we see that uniform convergence is again necessary and sufficient.

**Relationship to modern overparametrized models.** The inadequacies of the uniform convergence paradigm have been a topic of much study in modern deep learning theory (see e.g. [ZBH<sup>+</sup>21, NK19, BMR21, Bel21]). We may phrase the essential argument in the following way. Let  $\mathcal{C}$  be a certain “rich” concept class mapping  $X$  to  $\{0, 1\}$  (for concreteness). Let  $D$  be an unknown labeled distribution on  $X \rightarrow \{0, 1\}$ , and let  $S \sim D^m$  be a sample drawn from it. Let  $L$  and  $\widehat{L}_m$  denote the population and empirical 0-1 loss functionals, as before. Consider an ERM estimator  $\widehat{f}$  picked based on this sample:  $\widehat{f} \in \arg \min_{f \in \mathcal{C}} \widehat{L}_m(f)$ . We are interested in the generalization gap associated with  $\widehat{f}$ , namely the quantity  $|L(\widehat{f}) - \widehat{L}_m(\widehat{f})|$ . We would like to place a useful upper bound, say  $B$ , on this quantity.

The core observation is that certain classes  $\mathcal{C}$  (such as deep neural networks)

are rich enough that they can interpolate any sample of size  $m$ ; in this sense they are “overparametrized” relative to sample size  $m$ . In particular, they can fit even completely random labels. This of course means that  $\widehat{L}_m(\widehat{f}) = 0$  while  $L(\widehat{f}) = \frac{1}{2}$ . This in turn means that the bound  $B$  must be at least  $\frac{1}{2}$ . Note that this occurs without changing anything about the class  $\mathcal{C}$ , the marginal distribution  $D_X$  of  $D$ , or the training procedure (ERM). So any bound  $B$  that is purely a function of these quantities must be essentially vacuous; this includes uniform convergence bounds (e.g.,  $\sup_{f \in \mathcal{C}} |L(f) - \widehat{L}_m(f)| \leq R_m(\mathcal{C}, D_X)$ ) as well as algorithm-based bounds (e.g. those based on stability). Yet what is remarkable is that when the labels do satisfy some structure, e.g. when there exists  $f \in \mathcal{C}$  achieving error  $L(f) = \text{opt}(\mathcal{C}, D) < \frac{1}{2}$ , then we observe (provably or empirically) that the generalization gap is in fact relatively small, and  $\widehat{f}$  performs comparably with  $f$ . This phenomenon, sometimes referred to as “benign overfitting”, occurs not only with deep neural networks but in fact also already (in a sense) with linear regression [BLLT20, HMRT22]; we shall not attempt a summary of known results here but direct the reader to e.g. [BMR21, Bel21].

What the results in this chapter point out is that a version of this phenomenon also occurs in a strong, provable sense with classical examples such as convex sets in Gaussian space or monotone functions over the Boolean hypercube. These classes  $\mathcal{C}$  are also capable of interpolating a random sample of size  $m = 2^{\Omega(d)}$ ; and yet there exist estimators  $\widehat{f}$  that achieve error  $\text{opt}(\mathcal{C}, D) + \epsilon$  using sample complexity only  $2^{\widetilde{O}(\overline{d}/\text{poly}(\epsilon))}$  [BT96, KOS08]. These estimators are not based on ERM and do not lie strictly in  $\mathcal{C}$ ; instead, they are low-degree (specifically, degree- $O(\overline{d}/\text{poly}(\epsilon))$ ) polynomial approximators of functions in  $\mathcal{C}$  (as in Theorem 6.2.2). Such polynomial approximators essentially constitute a small (improper) cover of the class  $\mathcal{C}$  (w.r.t. the  $L^1(D_X)$  metric  $\rho(f, p) = \mathbb{E}_{D_X}[|f - p|]$ ). The improved sample complexity we obtain by such methods may be explained by the fact that to obtain generalization, it suffices to have uniform convergence over this cover as opposed to all of  $\mathcal{C}$  (as in the metric entropy characterization of [BI91]).

## 6.7.2 Implications for sandwiching degree

Another somewhat surprising consequence of the lower bounds in Section 6.6.2.1 is that the classes of convex sets over  $\mathcal{N}(0, I_d)$  and monotone functions over  $\text{Unif } \mathcal{F} \text{ } 1\mathcal{G}^d$  cannot admit sandwiching polynomials of degree  $o(d/\log d)$  and error even  $\epsilon = \Theta(1)$  unless they have very large coefficients. This is simply because any such sandwiching polynomials, if they have reasonable coefficients and if the distribution satisfies some concentration properties, will tend to allow the moment matching algorithm (Theorem 6.4.5) to succeed. More generally, we obtain the following surprising connection between Rademacher complexity and sandwiching degree as a direct corollary of Theorems 6.4.5 and 6.6.2.

**Corollary 6.7.1.** *Let  $\epsilon > 0$ , let  $D_X$  be a distribution on  $X$ , and let  $\mathcal{C}$  be a concept class mapping  $X$  to  $\mathcal{F} \text{ } 1\mathcal{G}$ . Let  $M$  be such that  $R_M(\mathcal{C}) \leq 5\epsilon$ , and assume  $M = \Theta(1/\epsilon^2)$ . Consider any degree and slack parameters  $k \geq 2 \in \mathbb{N}$ ,  $\Delta \geq \mathbb{R}_+^{j(k,d)}$  such that each  $f \in \mathcal{C}$  admits degree- $k$  sandwiching polynomials  $p_l \leq f \leq p_u$  satisfying*

$$\mathbb{E}_{D_X} [p_u - f] + h\Delta, |p_u| \leq \frac{\epsilon}{2}, \quad \mathbb{E}_{D_X} [f - p_l] + h\Delta, |p_l| \leq \frac{\epsilon}{2}.$$

*Let  $m$  be the sample complexity of testing with high probability whether the degree- $k$  empirical moments of  $D_X$  are within  $\Delta$  of their true moments. Then we must have  $m = \Omega(\overline{M}^{\frac{1}{k}})$ .*

Let us illustrate this in the cases where  $\mathcal{C}, D_X$  are either convex sets over  $\mathcal{N}(0, I_d)$  or monotone functions over  $\text{Unif } \mathcal{F} \text{ } 1\mathcal{G}^d$ . Consider any degree and slack parameters  $k, \Delta$ , and let  $\delta = \min_{I \subseteq [1, d], |I| = k} \Delta_I$ . In both cases, one can check with high probability whether the degree- $k$  empirical moments of  $D_X$  are within  $\Delta$  of their true moments using sample complexity at most  $m = d^{O(k)} \text{poly}(1/\delta)$  (for  $\text{Unif } \mathcal{F} \text{ } 1\mathcal{G}^d$  this is immediate by boundedness, while for  $\mathcal{N}(0, I_d)$  we appeal to Lemma 6.5.8). Now suppose that each  $f \in \mathcal{C}$  admitted degree- $k$   $(\epsilon/4)$ -sandwiching polynomials  $p_l \leq f \leq p_u$  satisfying  $\mathbb{E}_{D_X} [f - p_l], \mathbb{E}_{D_X} [p_u - f] \leq \epsilon/4$  and also such that their coefficients are bounded in magnitude by  $d^{O(k)}$ . Then clearly we can pick  $\Delta$  sufficiently small so that

$h\Delta, j p_{l j} i, h\Delta, j p_{u j} i \leq \epsilon/4$  while still ensuring  $\delta \leq \epsilon d^{-O(k)}$ . This means  $m$  as defined earlier is  $d^{O(k)} \text{poly}(1/\epsilon)$ . Thus for this choice of  $\Delta$ , both conditions (a) and (b) of Theorem 6.4.5 hold, and we obtain a testable learning algorithm with sample complexity  $m = d^{O(k)} \text{poly}(1/\epsilon)$ . For  $\epsilon = 0.1$ , say, we know by Section 6.6.2.1 that the required sample complexity for this task is  $2^{\Omega(d)}$ . Thus we see that  $k$  must necessarily be  $\Omega(d/\log d)$ .

The only way for sandwiching polynomials to exist despite this obstacle is by having unusually large coefficients (on the scale of  $d^{\omega(k)}$ ). Most reasonable approaches to constructing sandwiching polynomials will tend to ensure some boundedness of coefficients (indeed, this is true whenever such polynomials are constructed out of univariate polynomials that are bounded on a bounded domain, see e.g. [She12, Lemma 4.1]). Therefore, we regard this as good evidence in favor of a lower bound on the sandwiching degree for these classes.

# Appendix A: Superpolynomial Lower Bounds for Learning One-Layer Neural Networks using Gradient Descent

## A.1 Bounding the function norms under the Gaussian

Our goal in this section will be to give lower bounds on the norms of the functions in  $\mathcal{C}_{\text{orth}}(n, k)$ , which is a technical requirement for our results to hold (see Lemma 2.4.4 and Corollary 2.4.6). Note that when learning with respect to  $L_2$  error, such a lower bound is necessary if we wish to state SQ lower bounds, since if the target had small norm, say  $\|f\|_{k_D} \leq \epsilon$ , then the zero function trivially achieves  $L_2$  error  $\epsilon$ .

All inner products and norms in this section will be with respect to the standard Gaussian,  $\mathcal{N}(0, I)$ . Since we will fix  $S$  throughout, for our purposes the only relevant part of the input is  $x_S$  and so we drop the subscripts and let  $g = g_S, f = f_S$  and  $x = x_S$ , so that  $g$  and  $f$  are functions of  $x \in \mathbb{R}^k$ . Our approach will be as follows. In order to prove a norm lower bound on  $f$ , we will prove an anticoncentration result for  $g$ . To this end we first calculate the second moment of  $g$  in terms of the Hermite coefficients of  $\phi$ .

**Lemma A.1.1.** *Under the distribution  $\mathcal{N}(0, I_n)$ , let the Hermite representation of  $\phi$  be  $\phi(x) = \sum_{i=0}^1 \hat{\phi}_i \tilde{H}_i(x)$ , where  $\tilde{H}_i(x)$  is the  $i^{\text{th}}$  normalized probabilists' Hermite polynomial. Then*

$$\mathbb{E} [g(x)^2] = 4^k \sum_{i=0}^1 \frac{\hat{\phi}_i^2}{k^i} \sum_{\substack{i_1 + \dots + i_k = i \\ i_1, \dots, i_k \text{ are odd}}} \binom{i}{i_1, \dots, i_k}.$$

*Proof.* We use  $\mathbb{E}$  in this proof instead of  $\mathbb{E}_x \sim \mathcal{N}(0, I_n)$  for simplicity. Then we have

$$\begin{aligned}
& \mathbb{E}[g(x)^2] \\
&= \mathbb{E} \left[ \sum_{\alpha, \beta \in \mathcal{F}} \chi(\alpha) \phi \left( \frac{\alpha x_S}{k} \right) \right] \left[ \sum_{\beta \in \mathcal{F}} \chi(\beta) \phi \left( \frac{\beta x_S}{k} \right) \right] \\
&= \sum_{\alpha, \beta \in \mathcal{F}} \prod_{l=1}^k \alpha_l \beta_l \mathbb{E} \left[ \phi \left( \frac{\alpha x_S}{k} \right) \phi \left( \frac{\beta x_S}{k} \right) \right] \\
&= \sum_{\alpha, \beta \in \mathcal{F}} \prod_{l=1}^k \alpha_l \beta_l \mathbb{E} \left[ \sum_{i, j=0}^k \widehat{\phi}_i \widehat{\phi}_j \widetilde{H}_i \left( \frac{\alpha x_S}{k} \right) \widetilde{H}_j \left( \frac{\beta x_S}{k} \right) \right] \\
&= \sum_{\alpha, \beta \in \mathcal{F}} \prod_{l=1}^k \alpha_l \beta_l \sum_{i, j=0}^k \widehat{\phi}_i \widehat{\phi}_j \mathbb{E} \left[ \widetilde{H}_i \left( \frac{\alpha x_S}{k} \right) \widetilde{H}_j \left( \frac{\beta x_S}{k} \right) \right].
\end{aligned}$$

Since  $x \sim \mathcal{N}(0, I_k)$ ,  $\frac{h\alpha x_S^i}{k}$  and  $\frac{h\beta x_S^i}{k}$  are both standard Gaussian and have correlation  $\frac{h\alpha_i \beta_i}{k}$ , we then apply the following well-known property of the Hermite polynomials.

$$\mathbb{E}_{(a,b)^T \sim \mathcal{N}\left(0, \begin{pmatrix} 1 & \rho \\ \rho & 1 \end{pmatrix}\right)} \widetilde{H}_i(a) \widetilde{H}_j(b) = \delta_{i,j} \rho^i,$$

where  $\delta_{i,j}$  is the Dirac delta function.

$$\begin{aligned}
\mathbb{E}[g(x)^2] &= \sum_{\alpha, \beta \in \mathcal{F}} \prod_{l=1}^k \alpha_l \beta_l \sum_{i=0}^k \widehat{\phi}_i^2 \left( \frac{\alpha \beta}{k} \right)^i \\
&= \sum_{w, \theta \in \mathcal{F}} \prod_{l=1}^k w_l \sum_{i=0}^k \widehat{\phi}_i^2 \left( \frac{\sum_{l=1}^k w_l}{k} \right)^i \\
&= 2^k \sum_{w \in \mathcal{F}} \prod_{l=1}^k w_l \sum_{i=0}^k \widehat{\phi}_i^2 \left( \frac{\sum_{l=1}^k w_l}{k} \right)^i,
\end{aligned}$$

where  $w_i = \alpha_i \beta_i$  and  $\theta_i = w_i \alpha_i$ . Note that Assumption 2.3.3 implies that  $\sum_{i=0}^k \widehat{\phi}_i^2 <$



7 , the series above is absolute convergent. Then,

$$\begin{aligned}
& \mathbb{E} [g(x)^2] \\
&= 2^k \sum_{i=0}^{\infty} \frac{\widehat{\phi}_i^2}{k^i} \sum_{w \in \mathbb{Z}^k} \prod_{l=1}^k w_l \left( \frac{\sum_{l=1}^k w_l}{k} \right)^i \\
&= 2^k \sum_{i=0}^{\infty} \frac{\widehat{\phi}_i^2}{k^i} \sum_{w \in \mathbb{Z}^k} \prod_{l=1}^k w_l \sum_{i_1 + \dots + i_k = i} \prod_{l=1}^k w_l^{i_l} \binom{i}{i_1, \dots, i_k} \\
&= 2^k \sum_{i=0}^{\infty} \frac{\widehat{\phi}_i^2}{k^i} \sum_{i_1 + \dots + i_k = i} \binom{i}{i_1, \dots, i_k} \sum_{w \in \mathbb{Z}^k} \prod_{l=1}^k w_l^{i_l+1} \\
&= 2^k \sum_{i=0}^{\infty} \frac{\widehat{\phi}_i^2}{k^i} \sum_{i_1 + \dots + i_k = i} \binom{i}{i_1, \dots, i_k} \prod_{l=1}^k [1^{i_l+1} + (-1)^{i_l+1}] \\
&= 4^k \sum_{i=0}^{\infty} \frac{\widehat{\phi}_i^2}{k^i} \sum_{\substack{i_1 + \dots + i_k = i \\ i_1, \dots, i_k \text{ are odd}}} \binom{i}{i_1, \dots, i_k}
\end{aligned}$$

since we consider all distinct monomials in  $(\sum_{l=1}^k w_l)^i$ . Note that  $\sum_{i_1, \dots, i_k \text{ are odd}} \binom{i}{i_1, \dots, i_k}$  is always non-negative and is positive iff  $i = k$  and  $i = k \pmod{2}$ .  $\square$

### A.1.1 ReLU Activation

The goal of this section is to give a lower-bound of  $\|g\|$  for  $\phi = \text{ReLU}$  under the standard Gaussian distribution  $\mathcal{N}(0, I)$ . To this end, we prove an anti-concentration for  $g$ . We first give a lower bound on  $\|g\|$  based on the Hermite coefficients of  $\phi$ . If  $g$  were bounded, this alone would imply anti-concentration as in Section A.1.2. But since it is not, we first introduce  $g^T$ , where all activations are truncated at some  $T$ . We pick  $T$  large enough that  $g$  and  $g^T$  behave almost identically over  $\mathcal{N}(0, I)$ . We then show a lower bound on  $\|g^T\|$ , translate that into an anticoncentration result for  $g^T$ , and finally into one for  $g$ .

Let  $T > 0$  be some constant to be determined later. Let

$$\text{ReLU}^T(x) = \min(\text{ReLU}(x), T)$$

and

$$g^T(x) = \sum_{w \sim 1g^k} \chi(w) \text{ReLU}^T\left(\frac{xw}{k}\right).$$

The following lemma from [GKK19] describes the Hermite coefficients of ReLU.

**Lemma A.1.2.**

$$\text{ReLU}(x) = \sum_{i=0}^1 c_i \tilde{H}_i(x)$$

where

$$\begin{aligned} c_0 &= \sqrt{\frac{1}{2\pi}}, & c_1 &= \frac{1}{2}, \\ c_{2i-1} &= 0, & c_{2i} &= \frac{H_{2i}(0) + 2iH_{2i-2}(0)}{\sqrt{2\pi}(2i)!} \quad \text{for } i \geq 2. \end{aligned}$$

In particular,  $c_{2i}^2 = \Theta(i^{-2.5})$ .

We can now derive a lower bound on the norm of  $g$ .

**Lemma A.1.3.** *When  $k$  is even,*

$$\|g\| = \Omega\left(\left(\frac{4}{e}\right)^{\left(\frac{1}{2}+o(1)\right)k}\right).$$

*Proof.* Due to Lemma A.1.1,

$$\begin{aligned} \mathbb{E}[g(x)^2] &= 4^k \sum_{i=0}^k \frac{c_i^2}{k^i} \sum_{\substack{i_1+\dots+i_k=i \\ i_1, \dots, i_k, \text{ are odd}}} \binom{i}{i_1, \dots, i_k} \\ &= \frac{4^k c_k^2}{k^k} \sum_{\substack{i_1+\dots+i_k=k \\ i_1, \dots, i_k, \text{ are odd}}} \binom{k}{i_1, \dots, i_k} \\ &= \frac{4^k c_k^2 k!}{k^k}. \end{aligned}$$

The lemma then follows by the Stirling's approximation,

$$n! \sim \sqrt{2\pi n} \left(\frac{n}{e}\right)^n.$$

and the bound on the Hermite coefficients,

$$c_k^2 = \Theta(k^{-2.5}).$$

□

For the difference of  $g(x)$  and  $g^T(x)$ , we have

**Lemma A.1.4.**

$$\|g - g^T\| \leq 2^k e^{\frac{T^2}{4}} \sqrt{T^2 + 1} \frac{T}{2\pi}$$

*Proof.* Let  $\text{ReLU}_w(x)$  be shorthand for  $\text{ReLU}(\frac{x-w}{k})$ , and similarly  $\text{ReLU}_w^T$ . Observe that by the triangle inequality,

$$\begin{aligned} \|g - g^T\| &= \left\| \sum_{w \in \mathcal{W}} \chi(w) (\text{ReLU}_w - \text{ReLU}_w^T) \right\| \\ &\leq \sum_{w \in \mathcal{W}} \|\text{ReLU}_w - \text{ReLU}_w^T\| \\ &= 2^k \|\text{ReLU} - \text{ReLU}^T\|_{\mathcal{N}(0,1)}, \end{aligned}$$

where the last equality holds because for any unit vector  $v$  and  $x \sim \mathcal{N}(0, I)$ ,  $x \cdot v$  has the distribution  $\mathcal{N}(0, 1)$ . Now,

$$\|\text{ReLU} - \text{ReLU}^T\|_{\mathcal{N}(0,1)}^2 = \int_T^1 (x - T)^2 p(x) dx,$$

where  $p(x)$  is the probability density function of  $\mathcal{N}(0, 1)$ . Note that  $p'(x) = -xp(x)$ .

We have

$$\begin{aligned} \int_T^1 x^2 p(x) dx &= \int_T^1 x d(p(x)) \\ &= x p(x) \Big|_T^1 + \int_T^1 p(x) dx && \text{(integration by parts)} \\ &= T p(T) + \mathbb{P}_{x \sim \mathcal{N}(0,1)}(x > T), \\ \int_T^1 x p(x) dx &= p(x) \Big|_T^1 = p(T), \\ \int_T^1 p(x) dx &= \mathbb{P}_{x \sim \mathcal{N}(0,1)}(x > T) = e^{-\frac{T^2}{2}}. \end{aligned}$$

Thus,

$$\begin{aligned} & \mathbb{E} [g(x) - g^T(x)]^2 \\ & 4^k \left[ (T^2 + 1) \mathbb{P}_{x \sim N(0,1)}(x > T) + T p(T) \right] \\ & 4^k e^{-\frac{T^2}{2}} \left( T^2 + 1 - \frac{T}{2\pi} \right). \end{aligned}$$

□

**Lemma A.1.5.**

$$\mathbb{P}[g(x) \neq g^T(x)] \leq 2^k e^{-\frac{T^2}{2}}.$$

*Proof.* For any  $w \geq 1$  and  $g^k$ ,

$$\begin{aligned} & \mathbb{P}_{x \sim N(0,1)} \left[ \text{ReLU}\left(\frac{xw}{k}\right) \neq \text{ReLU}^T\left(\frac{xw}{k}\right) \right] \\ & = \mathbb{P}_{t \sim N(0,1)} [t > T] \\ & \leq e^{-\frac{T^2}{2}}. \end{aligned}$$

The lemma follows by a union bound. □

**Lemma A.1.6.**

$$\mathbb{P} [ \|g(x)\| \geq 1 ] = \Omega(\exp(-\Theta(k))).$$

*Proof.* For large enough  $T = \Omega(k)$ , it holds from Lemmas A.1.3 and A.1.4 that

$$\|g^T\| = \Omega\left(\left(\frac{4}{e}\right)^{\left(\frac{1}{2} + o(1)\right)k}\right).$$

Since  $|g^T(x)| \leq T 2^k$ ,

$$\|g^T\|^2 = \mathbb{E}[g^T(x)^2] \leq 1 + \mathbb{P}[|g^T(x)| \geq 1] + (T 2^k)^2 \mathbb{P}[|g^T(x)| \geq 1],$$

so that

$$\mathbb{P} [ \|g^T(x)\| \geq 1 ] \leq \frac{\Omega\left(\left(\frac{4}{e}\right)^{(1+o(1))k}\right)}{(T 2^k)^2} = \Omega(\exp(-\Theta(k))) \quad (\text{A.1.1})$$

Using Eq. (A.1.1) with Lemma A.1.5,

$$\mathbb{P}[|g(x)| \geq 1] = \Omega(\exp(-\Theta(k)))$$

for large enough  $T = \Omega(k)$ . □

The lower bound on  $\|kf\|$  now follows easily.

**Corollary A.1.7.**

$$\|kf\| = \Omega(\exp(-\Theta(k))).$$

*Proof.* Since  $f = \psi \circ g$ , from Lemma A.1.6 and the fact that  $\psi$  is odd and increasing, we have that

$$\begin{aligned} \|kf\| &= |\psi(1)| \mathbb{P}[g(x) \geq 1] + |\psi(-1)| \mathbb{P}[g(x) \leq -1] \\ &= \psi(1) \mathbb{P}[|g(x)| \geq 1] \\ &= \Omega(\exp(-\Theta(k))). \end{aligned}$$

□

### A.1.2 Sigmoid Activation

Here we consider  $g$  and  $f$  with  $\phi(x) = \sigma(x) = \frac{1}{1+e^{-x}}$ . For the asymptotic bound of Hermite polynomial coefficients, we need the following theorem from [Boy84].

**Theorem A.1.8.** *For a function  $f(z)$  whose convergence is limited by simple poles at the roots of  $z^2 = -\gamma^2$  with residue  $R$ , the non-zero expansion coefficients  $\hat{f}_{a_n}$  of  $f(z)$  as a series of normalized Hermite functions have magnitudes asymptotically given by*

$$|\hat{f}_{a_n}| \sim 2^{\frac{5}{4}} \pi^{\frac{1}{2}} |R| n^{-\frac{1}{4}} e^{-\gamma(2n+1)^{\frac{1}{2}}},$$

Here the normalized Hermite function  $\hat{f}_{\psi_n}(x)g_{n,2\mathbb{N}}$  is defined by

$$\psi_n(z) = e^{-\frac{z^2}{2}} \pi^{-\frac{1}{4}} \tilde{H}_n\left(\frac{z}{\sqrt{2}}\right).$$

Applying this to  $f(x) = e^{-\frac{x^2}{2}} \sigma(\sqrt{2}x)$  and translating the Hermite coefficients for the series in terms of Hermite functions to those in terms of Hermite polynomials, we have

**Lemma A.1.9.**

$$\sigma(x) = \sum_{i=0}^{\infty} c_i \tilde{H}_i(x),$$

where  $c_0 = 0.5, c_{2i} = 0$  for  $i \geq 1$  and all non-zero odd terms satisfies

$$c_{2i-1} = e^{-\binom{p}{i}}.$$

**Corollary A.1.10.** *There is an infinite increasing sequence  $\{k_i\}_{i \in \mathbb{N}}$  such that  $k_i$ 's are all odd and*

$$c_{k_i} = e^{-\binom{p}{k_i}}.$$

*Proof.* It follows simply from the fact that  $\sigma$  is not a polynomial and there should be infinitely many non-zero terms in  $\sum_{k \in \mathbb{N}} c_k g_{k,2\mathbb{N}}$ .  $\square$

*Remark A.1.11.* Experimental evidence strongly indicates that in fact all odd Hermite coefficients of sigmoid are nonzero and decay as above, but this is laborious to formally establish. So we state our norm lower bound only for  $k \geq k_i$  (and the associated  $n \geq 2^{k_i}$ , since we end up taking  $k = \log n$ ). Since this is nevertheless an infinite sequence, it still establishes that no better asymptotic bound holds.

Similar to Lemma A.1.3, we can derive a lower bound of  $\|g_k\|$  for some  $k$ 's.

**Lemma A.1.12.** *For  $k \geq k_i$ ,*

$$\|g(x)_k\| = \Omega \left( \left( \frac{4}{e} \right)^{\binom{1}{2} + o(1)} k \right).$$

*Proof.* Due to Lemma A.1.1,

$$\begin{aligned} \mathbb{E} [g(x)^2] &= 4^k \sum_{i=0}^k \frac{c_i^2}{k^i} \sum_{\substack{i_1 + \dots + i_k = i \\ i_1, \dots, i_k, \text{ are odd}}} \binom{i}{i_1 \dots i_k} \\ &= \frac{4^k c_k^2}{k^k} \sum_{\substack{i_1 + \dots + i_k = k \\ i_1, \dots, i_k, \text{ are odd}}} \binom{k}{i_1 \dots i_k} \\ &= \frac{4^k c_k^2 k!}{k^k}. \end{aligned}$$

Using Stirling's approximation,

$$n! \sim \sqrt{2\pi n} \left(\frac{n}{e}\right)^n,$$

and Corollary A.1.10,

$$c_k = e^{-\binom{k}{2}},$$

we obtain

$$\mathbb{E} [g(x)^2] = \Omega \left( \frac{4^k e^{-\binom{k}{2}}}{k^k} \left(\frac{k}{e}\right)^k e^{-\binom{k}{2}} \right)$$

and hence

$$\mathbb{E} [g(x)^2] = \Omega \left( \left(\frac{4}{e}\right)^{(1+o(1))k} \right).$$

□

**Lemma A.1.13.** For  $k \geq 2$  and  $g_i \in \mathcal{G}_{2N}$ ,

$$\mathbb{P}(|g(x) - 1| \geq 1) = \Omega(\exp(-\Theta(k))).$$

*Proof.* Since  $|g(x) - 1| \leq 2^k$ ,

$$k g^2 = \mathbb{E}[g(x)^2] - \mathbb{P}[|g(x) - 1| \geq 1] + (2^k)^2 \mathbb{P}[|g(x) - 1| \geq 1],$$

and so

$$\mathbb{P}(|g(x) - 1| \geq 1) = \frac{\Omega\left(\left(\frac{4}{e}\right)^{(1+o(1))k}\right) - 1}{(2^k)^2}.$$

The lemma then follows. □

Using the same argument as Corollary A.1.7, we have the following bound.

**Corollary A.1.14.**

$$\|fk\| = \Omega(\exp(-\Theta(k))).$$

### A.1.3 General activations

It is not hard to see that the norm analysis of ReLU and sigmoid extends to any activation function for which a suitable lower bound on the Hermite coefficients holds, and which is either bounded or grows at a polynomial rate, so that under the standard Gaussian it behaves essentially identically to its truncated form. In particular, a lower bound of  $\alpha^{-j}$  for any constant  $\alpha < 4/e$  on the  $j^{\text{th}}$  Hermite coefficient suffices to give  $\|fk\| \geq \exp(-\Theta(k))$ , by the same argument as in Lemma A.1.3 and Lemma A.1.12. This then suffices to give  $\|fk\| \geq \exp(-\Theta(k))$ , as above.

In fact, even a very weak lower bound on  $\|fk\|$  yields *some* superpolynomial bound on learning. Suppose we only had  $\|fk\| \geq 1/\exp(\exp(\Theta(k)))$ , for instance. Then we can take  $k = \log \log n$  and have  $\|fk\| \geq 1/\text{poly}(n)$  and still obtain a lower bound of  $n^{\log \log n} = n^{\omega(1)}$  (see Theorem 2.3.9). Any lower bound on  $\|fk\|$  will be a function only of  $k$ , so a similar argument applies.

## A.2 SQ lower bound for real-valued functions proof

We give a self-contained variant of the elegant proof of [Szö09] for the reader's convenience. For simplicity, we include the 0 function in our class  $\mathcal{C}$  — this can only negligibly change the SDA, and it makes the core argument cleaner.

**Theorem A.2.1.** *Let  $D$  be a distribution on  $X$ , and let  $\mathcal{C}$  be a real-valued concept class over a domain  $X$  such that  $0 \in \mathcal{C}$ , and  $\|ck\|_D > \epsilon$  for all  $c \in \mathcal{C}, c \neq 0$ . Consider any SQ learner that is allowed to make only inner product queries to an SQ oracle for the labeled distribution  $D_c$  for some unknown  $c \in \mathcal{C}$ . Let  $d = \text{SDA}_D(\mathcal{C}, \gamma)$ . Then any such SQ learner needs at least  $d/2$  queries of tolerance  $\frac{\epsilon}{\gamma}$  to learn  $\mathcal{C}$  up to  $L_2$*



error  $\epsilon$ .

*Proof.* Consider the adversarial strategy where we respond to every query  $h : X \rightarrow \mathbb{R}$  ( $\|h\|_D = 1$ ) with 0. This corresponds to the true expectation if the target were the 0 function. By the norm lower bound, outputting any other  $c$  would then mean  $L_2$  error greater than  $\epsilon$ . Thus we must rule out all other  $c \in \mathcal{C}$ .

Let  $\tau = \frac{\rho_D}{\gamma}$ . If  $h_k$  is the  $k^{\text{th}}$  query, let  $S_k = \{c \in \mathcal{C} \mid \langle h_k, c \rangle_D > \tau\}$  be the functions ruled out by our response of 0. (A similar argument will hold for  $S_k^0 = \{c \in \mathcal{C} \mid \langle h_k, c \rangle_D < -\tau\}$ .) Let  $\Phi = \langle h_k, \sum_{c \in S_k} c \rangle_D$ . We claim that  $\| \sum_{c \in S_k} c \|_D \geq \tau |S_k|/d$ . Suppose not. Then  $\rho_D(S_k) < \tau$  by Definition 2.2.4, and

$$\begin{aligned} \Phi &= \langle h_k, \sum_{c \in S_k} c \rangle_D \\ &= \sqrt{\sum_{c \in S_k} \langle h_k, c \rangle_D^2} \\ &= \sqrt{|S_k| \rho_D(S_k)} \\ &< \frac{\rho_D}{\gamma} |S_k|, \end{aligned}$$

contradicting the fact that  $\Phi \geq \tau |S_k|$  by definition of  $S_k$ .

Similarly  $|S_k^0| \leq \frac{\rho_D}{\gamma} |S_k^0|$ . Thus we rule out at most a  $2/d$  fraction of functions with each query, and hence need at least  $d/2$  queries to rule out all other possibilities.  $\square$

# Appendix B: Hardness of Noise-Free Learning for Two-Hidden-Layer Neural Networks

## B.1 Barriers for constructing $N_3$

We briefly discuss why one natural approach to constructing  $N_3$  satisfying the ideal properties in Eq. (3.3.2) ultimately requires two hidden layers rather than one, unlike the construction we give in Sections 3.3.2 and 3.3.3.

The most straightforward way to ensure that a function of  $s_1, \dots, s_d, t$  vanishes whenever there exists  $j$  for which  $s_j = 1$  would be to threshold on  $\sum s_j$ , e.g. by taking  $\text{ReLU}(1 - \sum_j s_j)$ . While this function is a one-hidden-layer ReLU network, it is unclear how to modify it to satisfy the remaining desiderata in (3.3.2) while preserving the fact that it has only one hidden layer. We note that [DV21] takes this approach of thresholding on  $\sum_j s_j$  but uses two hidden layers.

Here we informally argue that such an approach inherently requires an extra hidden layer. That is, we argue that no function  $N : \mathbb{R}^2 \rightarrow \mathbb{R}$  that takes as inputs  $s, \sum_j s_j$  and  $t$  and satisfies (3.3.2) can be implemented as a one-hidden-layer network. Concretely,  $N(s, t)$  must vanish whenever  $s = 1$  or  $t \geq 1/g$ . Any function computed by a one-hidden-layer ReLU network of the form  $(s, t) \mapsto \sum_i \text{ReLU}(a_i s + b_i t - c_i)$ , unless if it is affine linear, must in general be nowhere smooth (i.e. have a discontinuous gradient) along the entire line where a particular neuron of the network vanishes. In our example, these are the lines  $f(s, t) : a_i s + b_i t = c_i/g$ . But this means that such a line cannot intersect the region  $f(s, t) : s = 1/g$ , as otherwise it would be zero (hence smooth) on an infinite segment of the line. This can only happen if  $b_i = 0$ , i.e. none of the neurons of  $N$  depend on  $t$ . Such a network clearly cannot satisfy (3.3.2).

## B.2 Supporting lemmas for Section 3.3

**Lemma B.2.1.** For any  $0 \leq S < m \leq d$ ,

$$\sum_{i=S}^m (-1)^{m-i} \binom{d-i-2}{d-m-2} \binom{d-m-1}{i-S} = 0. \quad (\text{B.2.1})$$

*Proof.* We will show that for any integers  $j \leq \ell - 1$ ,

$$\sum_{k=0}^{\ell} (-1)^k \binom{j-k}{\ell-1} \binom{\ell}{k} = 0. \quad (\text{B.2.2})$$

We would like to substitute  $\ell = d - m - 1$  and  $j = d - 2 - S$ . Note that this is valid as we can assume without loss of generality that  $d - m - 1 \geq 1$  (otherwise  $\binom{d-m-1}{i-S} = 0$  on the right-hand side of (B.2.1)), and  $j \leq \ell$  by our assumption that  $S < m$ . We conclude the identity

$$0 = \sum_{k=0}^{d-m-1} (-1)^k \binom{d-2-S-k}{d-m-2} \binom{d-m-1}{k} = \sum_{i=S}^{d-m-1+S} (-1)^{i-S} \binom{d-i-2}{d-m-2} \binom{d-m-1}{i-S}, \quad (\text{B.2.3})$$

where the second step is by the change of variable  $i = k + S$ . If  $d - m - 1 + S \leq m$ , then note that all summands  $m < i \leq d - m - 1 + S$  vanish because in that case  $d - i - 2 < d - m - 2$  and so  $\binom{d-i-2}{d-m-2} = 0$ . If  $d - m - 1 + S > m$ , then note that all summands  $d - m - 1 + S < i \leq m$  vanish because in that case  $d - m - 1 < i - S$  and so  $\binom{d-m-1}{i-S} = 0$ . We conclude that (B.2.3) is equal, up to a sign, to the left-hand side of (B.2.1), so we'd be done.

It remains to establish (B.2.2), which we do by following an argument due to [Ear19]. Observe that the left-hand side of (B.2.2) is simply counting via inclusion-exclusion the number of subsets of  $\{1, \dots, j\}$  of size  $\ell - 1$  which contain  $1, \dots, \ell$ . Indeed, the  $k = 0$  summand counts all subsets of size  $\ell - 1$ . The  $k = 1$  summands subtract out the contribution, for every  $1 \leq x \leq \ell$ , from the subsets of size  $\ell - 1$  which contain  $x$ . The  $k = 2$  summands add back the contribution, for every distinct  $1 \leq x < y \leq \ell$ , from the subsets of size  $\ell - 1$  which contain both of  $x, y$ , etc.  $\square$

**Lemma B.2.2.** For any integers  $m \geq 3$  and  $a \in \{0, 1, 2\}$ ,

$$\sum_{i=1}^m \sum_{j=1}^{m+1-i} (-1)^{m-i} \binom{d-i}{m-i} \binom{j-1}{j+1} \binom{d-1}{i-1} j^a = \mathbb{1}[a=0] \quad (\text{B.2.4})$$

$$\sum_{i=0}^m \sum_{j=1}^{m+1-i} (-1)^{m-i} \binom{d-i}{m-i} \binom{j-1}{j+1} \binom{d-1}{i} j^a = 0 \quad (\text{B.2.5})$$

*Proof.* By taking  $\ell = i + j$ , we can rewrite these sums as

$$S_{a,m} = \sum_{\ell=2}^{m+1} \sum_{i=1}^{\ell-1} (-1)^{m-i} \binom{d-1}{m+1} \binom{\ell}{\ell} \binom{d-1}{i-1} (\ell-i)^a \quad (\text{B.2.6})$$

$$T_{a,m} = \sum_{\ell=1}^{m+1} \sum_{i=0}^{\ell-1} (-1)^{m-i} \binom{d-1}{m+1} \binom{\ell}{\ell} \binom{d-1}{i} (\ell-i)^a \quad (\text{B.2.7})$$

We proceed by induction on  $m$ . The base cases follow from a direct calculation. By the change of variable  $\ell^\theta = \ell - 1$ , we can rewrite  $S_{a,m+1}$  as

$$\sum_{\ell^\theta=1}^{m+1} \sum_{i=1}^{\ell^\theta} (-1)^{m-i} \binom{d-1}{m+1} \binom{\ell^\theta}{\ell^\theta} \binom{d-1}{i-1} (\ell^\theta+1-i)^a \quad (\text{B.2.8})$$

$$= \sum_{\ell^\theta=1}^{m+1} \sum_{i=1}^{\ell^\theta} (-1)^{m-i} \binom{d-1}{m+1} \binom{\ell^\theta}{\ell^\theta} \binom{d-1}{i-1} (\ell^\theta-i)^a \quad (\text{B.2.9})$$

$$\sum_{\ell^\theta=1}^{m+1} \sum_{i=1}^{\ell^\theta} (-1)^{m-i} \binom{d-1}{m+1} \binom{\ell^\theta}{\ell^\theta} \binom{d-1}{i-1} \sum_{b=0}^{a-1} \binom{a}{b} (\ell^\theta-i)^b \quad (\text{B.2.10})$$

Note that the first term on the right-hand side differs from  $S_{a,m}$  only in the summands given by  $1 - i = \ell^\theta - m + 1$ , and those summands clearly vanish. We conclude that the first term on the right-hand side of (B.2.10) is exactly  $S_{a,m}$ . For the second term on the right-hand side of (B.2.10), the part coming from any  $0 < b - a - 1$  is also zero, so we thus get

$$= S_{a,m} \sum_{\ell^\theta=1}^{m+1} \sum_{i=1}^{\ell^\theta} (-1)^{m-i} \binom{d-1}{m+1} \binom{\ell^\theta}{\ell^\theta} \binom{d-1}{i-1} \quad (\text{B.2.11})$$

$$= S_{a,m} S_{0,m} \sum_{\ell^\theta=1}^{m+1} (-1)^{m-\ell^\theta} \binom{d-1}{m+1} \binom{\ell^\theta}{\ell^\theta} \binom{d-1}{\ell^\theta-1} \quad (\text{B.2.12})$$

$$= S_{a,m} (-1)^{m-1} \sum_{\ell^\theta=1}^{m+1} (-1)^{m-\ell^\theta} \binom{d-1}{m+1} \binom{\ell^\theta}{\ell^\theta} \binom{d-1}{\ell^\theta-1} \quad (\text{B.2.13})$$

$$= S_{a,m} = \mathbb{1}[a=0], \quad (\text{B.2.14})$$

where the penultimate step follows e.g. by applying the identity in [PSP17]. This completes the induction for  $S_{a,m}$ .

For  $T_{a,m}$ , note that by the change of variable  $i^\theta = i + 1$ ,

$$T_{a,m} = \sum_{\ell=1}^{m+1} \sum_{i^\theta=1}^{\ell} (-1)^{m-i^\theta} \binom{d-1}{m+1} \binom{\ell}{\ell} \binom{d-1}{i^\theta-1} (\ell - i^\theta + 1)^a \quad (\text{B.2.15})$$

$$= \sum_{\ell=2}^{m+1} \sum_{i^\theta=1}^{\ell-1} (-1)^{m-i^\theta} \binom{d-1}{m+1} \binom{\ell}{\ell} \binom{d-1}{i^\theta-1} (\ell - i^\theta + 1)^a + \sum_{\ell=1}^{m+1} (-1)^{m-\ell} \binom{d-1}{m+1} \binom{\ell}{\ell} \binom{d-1}{\ell-1} \quad (\text{B.2.16})$$

$$= \sum_{b=0}^a \binom{a}{b} S_{b,m} + 1 = 0, \quad (\text{B.2.17})$$

where in the second step we pulled out the summands corresponding to  $i^\theta = \ell$ , in the third step we used (B.2.14), and in the last step we used that for  $m \geq 3$ ,  $S_{b,m} = \mathbb{1}[b \notin [0, m-2]]$  for  $0 \leq b \leq 2$ .  $\square$

### B.3 SQ lower bound for the LWR functions

Here we prove an SQ lower bound for the LWR functions (Theorem 3.4.5) using a general formulation in terms of pairwise independent function families. To

our knowledge, this particular formulation has not appeared explicitly before in the literature, and was communicated to us by Bogdanov [Bog21]. A variant of this argument may be found in [BR17, §7.7].

**Definition B.3.1.** Let  $\mathcal{C}$  be a function family mapping  $X$  to  $Y$ , and let  $D$  be a distribution on  $X$ . We call  $\mathcal{C}$  an  $(1 - \eta)$ -pairwise independent function family if with probability  $1 - \eta$  over the choice of  $x, x^\theta$  drawn independently from  $D$ , the distribution of  $(f(x), f(x^\theta))$  for  $f$  drawn uniformly at random from  $\mathcal{C}$  is the product distribution  $\text{Unif}(Y) \times \text{Unif}(Y)$ .

**Lemma B.3.2.** Fix security parameter  $n$  and moduli  $p, q$ . The  $LWR_{n,p,q}$  function class  $\mathcal{C}_{LWR} = \{f_w : \mathbb{Z}_q^n \rightarrow \mathbb{Z}_p\}$  is  $(1 - \frac{2}{q^{n-1}})$ -pairwise independent with respect to  $\text{Unif}(\mathbb{Z}_q^n)$ .

*Proof.* This follows from the simple observation that whenever  $x, x^\theta \in \mathbb{Z}_q^n$  are linearly independent, the pair  $(w \cdot x \bmod q, w \cdot x^\theta \bmod q)$  for  $w \in \text{Unif}(\mathbb{Z}_q^n)$  is distributed as  $\text{Unif}(\mathbb{Z}_q) \times \text{Unif}(\mathbb{Z}_q)$ . For such  $x, x^\theta$ ,  $(f_w(x), f_w(x^\theta)) = (\frac{1}{p}bw \cdot x \bmod qe_p, \frac{1}{p}bw \cdot x^\theta \bmod qe_p)$  for  $f_w \in \text{Unif}(\mathcal{C}_{LWR})$  is distributed as  $\text{Unif}(\mathbb{Z}_p/p) \times \text{Unif}(\mathbb{Z}_p/p)$ . The probability that  $x, x^\theta \in \text{Unif}(\mathbb{Z}_q^n)$  are linearly dependent is at most

$$\mathbb{P}[x = 0] + \mathbb{P}[x \neq 0] \mathbb{P}[x^\theta \text{ is a multiple of } x] \leq \frac{1}{q^n} + \frac{q}{q^n} = \frac{2}{q^{n-1}}.$$

□

We can now prove full SQ lower bounds for any  $(1 - \eta)$ -pairwise independent function family as follows.

**Lemma B.3.3.** Let  $\mathcal{C}$  mapping  $X$  to  $Y$  be a  $(1 - \eta)$ -pairwise independent function family w.r.t. a distribution  $D$  on  $X$ . Let  $\phi : X \rightarrow Y \in [-1, 1]$  be any bounded query function. Then

$$\text{Var}_f \mathbb{E}_{x \sim D} [\phi(x, f(x))] \geq 2\eta.$$

*Proof.* Denote  $\mathbb{E}_{x \sim D}[\phi(x, f(x))]$  by  $\phi[f]$ . By some algebraic manipulations (with all subscripts denoting independent draws),

$$\mathbb{E}_f \text{Var}_{\text{Unif}(C)}[\phi[f]] = \mathbb{E}_f [\phi[f]^2] - (\mathbb{E}_f [\phi[f]])^2 \quad (\text{B.3.1})$$

$$= \mathbb{E}_f [\phi[f]\phi[f]] - \mathbb{E}_f [\phi[f]] \mathbb{E}_{f^\theta} [\phi[f^\theta]] \quad (\text{B.3.2})$$

$$= \mathbb{E}_{f, f^\theta} \left[ \mathbb{E}_x [\phi(x, f(x))] \mathbb{E}_{x^\theta} [\phi(x^\theta, f(x^\theta))] - \mathbb{E}_x [\phi(x, f(x))] \mathbb{E}_{x^\theta} [\phi(x^\theta, f^\theta(x^\theta))] \right] \quad (\text{B.3.3})$$

$$= \mathbb{E}_{x, x^\theta} \mathbb{E}_{f, f^\theta} [\phi(x, f(x))\phi(x^\theta, f(x^\theta)) - \phi(x, f(x))\phi(x^\theta, f^\theta(x^\theta))]. \quad (\text{B.3.4})$$

By  $(1 - \eta)$ -pairwise independence of  $C$ , the inner expectation vanishes with probability  $1 - \eta$  over the choice of  $x, x^\theta \sim D$ , and is at most  $2\eta$  otherwise. This gives the claim.  $\square$

**Theorem B.3.4.** *Let  $C$  mapping  $X$  to  $Y$  be a  $(1 - \eta)$ -pairwise independent function family w.r.t. a distribution  $D$  on  $X$ . For any  $f \in C$ , let  $D_f$  denote the distribution of  $(x, f(x))$  where  $x \sim D$ . Let  $D_{\text{Unif}(C)}$  denote the distribution of  $(x, y)$  where  $x \sim D$  and  $y = f(x)$  for  $f \sim \text{Unif}(C)$  (this can be thought of as essentially  $D \times \text{Unif}(Y)$ ). Any SQ learner able to distinguish the labeled distribution  $D_f$  for an unknown  $f \in C$  from the randomly labeled distribution  $D_{\text{Unif}(C)}$  using bounded queries of tolerance  $\tau$  requires at least  $\frac{\tau^2}{2\eta}$  such queries.*

*Proof.* Let  $\phi : X \times Y \rightarrow [-1, 1]$  be any query made by the learner. For any  $f \in C$ , let  $\phi[f]$  denote  $\mathbb{E}_{x \sim D}[\phi(x, f(x))] = \mathbb{E}_{(x, y) \sim D_f}[\phi(x, y)]$ . Consider the adversarial strategy where the SQ oracle responds to this query with  $\bar{\phi} = \mathbb{E}_{f \sim \text{Unif}(C)} \phi[f] = \mathbb{E}_{(x, y) \sim D_{\text{Unif}(C)}}[\phi(x, y)]$ . By Chebyshev's inequality and Lemma B.3.3,

$$\mathbb{P}_f \left[ |\phi[f] - \bar{\phi}| > \tau \right] \leq \frac{\text{Var}_f \text{Unif}(C) [\phi[f]]}{\tau^2} \leq \frac{2\eta}{\tau^2}.$$

So each such query only allows the learner to rule out at most a  $\frac{2\eta}{\tau^2}$  fraction of  $C$ . Thus to distinguish  $D_f$  from  $D_{\text{Unif}(C)}$ , the learner requires at least  $\frac{\tau^2}{2\eta}$  queries.  $\square$

Theorem 3.4.5 now follows easily as a corollary.

*Proof of Theorem 3.4.5.* It is not hard to see that learning  $\mathcal{C}_{\text{LWR}}$  up to squared loss  $1/16$  certainly suffices to solve the distinguishing problem in Theorem B.3.4. The claim now follows by Lemma B.3.2.  $\square$

*Remark B.3.5.* We remark that the argument in this section, specialized to the  $q = 2$  case, recovers the traditional SQ lower bound for parities (Theorem 3.4.3) without appealing to any notion of statistical dimension.



# Appendix C: Statistical-Query Lower Bounds via Functional Gradients

## C.1 SQ lower bound subtleties

### C.1.1 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  $\tau$ , the final desired  $L^2$  error  $\|f - f^*\|$  is  $\epsilon$  (which corresponds to  $L(f) - L(f^*) = \epsilon^2$ ; see Eq. (4.3.1)), and that the functions in  $\mathcal{C}$  satisfy  $\|f\| \leq \beta$  for all  $f \in \mathcal{C}$ . Then

1. We must have  $\tau < \epsilon$ . To see why, first note that for any query  $\phi$  and two functions  $f, g \in \mathcal{C}$ , a calculation shows that  $|E_{D_f}[\phi] - E_{D_g}[\phi]| = |\int \phi(x) (f(x) - g(x)) dx| \leq \beta \|f - g\|$ , 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  $\tau$  can tell them apart, but  $f$  is not  $\epsilon$ -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  $\beta$  since the functions must remain Boolean  $p$ -concepts. Rescaling would also increase the description length of the functions.

The lower bound in Theorem 4.2.2 (from [GGJ<sup>+</sup>20]) is proved by reducing a distinguishing problem to a learning problem. For technical reasons, we end up requiring  $\tau = \epsilon^2$ ,  $\epsilon = \beta/3$  for this reduction to go through. The points above show that these requirements are essentially necessary.

### C.1.2 The dependence of the query lower bound on the error $\epsilon$ and the tolerance $\tau$

The relationship between our query lower bounds, the desired error  $\epsilon$ , and the tolerance  $\tau$  may seem a little unusual at first sight, especially the fact that the lower bounds seem to grow weaker as  $\epsilon$  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  $D_0 = D \text{ Unif } 1g$ , using queries of tolerance at least  $\tau$ .

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  $\tau$ . In particular, using queries of tolerance  $\tau$ , we need at least  $\text{SDA}(\mathcal{C}, \tau^2)$  queries. This bound increases with  $\tau$ ; in fact it often scales as  $j\mathcal{C}j\tau^2$  (see ([GGJ+20, Theorem 4.5 and Lemma 2.6])).
2. The problem of learning up  $\mathcal{C}$  to error  $\epsilon$  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  $\tau$ , the query lower bound does not technically grow as a function of the error  $\epsilon$ : it applies uniformly for all  $\epsilon$  small enough that it allows the learner to solve the distinguishing problem. In other words,

there is a certain “threshold”  $\epsilon_0$  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  $\beta$  is such that  $kck \leq \beta$  for all  $c \geq \mathcal{C}$ .

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  $\tau \leq \epsilon^2$ ). If  $\tau \geq \epsilon$ , learning up to error  $\epsilon$  is simply impossible.

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

## C.2 Bounding the function norms of the [DKKZ20] construction

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

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

$$f(x) = \psi(g(x)),$$

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  $\mathbb{R}^n \rightarrow \mathbb{R}$  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<sup>+</sup>20].

**Lemma C.2.1.** Let  $g : \mathbb{R}^2 \rightarrow \mathbb{R}$  be as defined above, and assume  $m$  is even. Assume the standard Hermite expansion of  $\phi$  is given by  $\phi = \sum_a \widehat{\phi}_a H_a$ , where the  $H_a$  are the normalized probabilists' Hermite polynomials. Under  $N(0, I_2)$ ,

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

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

*Proof.* We have

$$\begin{aligned} kgk^2 = \mathbb{E}[g(x)^2] &= \frac{1}{4m^2} \sum_{i,j=1}^{2m} \binom{2m}{i} \binom{2m}{j} \mathbb{E}[\phi(x - w_i)\phi(x - w_j)] \\ &= \frac{1}{4m^2} \sum_{i,j=1}^{2m} \binom{2m}{i} \binom{2m}{j} \mathbb{E} \left[ \left( \sum_a \widehat{\phi}_a H_a(x - w_i) \right) \left( \sum_b \widehat{\phi}_b H_b(x - w_j) \right) \right] \\ &= \frac{1}{4m^2} \sum_{i,j=1}^{2m} \binom{2m}{i} \binom{2m}{j} \left( \sum_a \widehat{\phi}_a^2 \mathbb{E}[H_a(x - w_i)H_a(x - w_j)] \right). \end{aligned}$$

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  $x - w_i$  and  $x - w_j$  are both  $N(0, 1)$  with covariance  $\cos \frac{(i-j)\pi}{m}$ . Thus

$$\begin{aligned} kgk^2 &= \frac{1}{4m^2} \sum_{i,j=1}^{2m} \binom{2m}{i+j} \left( \sum_a \widehat{\phi}_a^2 \cos^a \frac{(i-j)\pi}{m} \right) \\ &= \frac{1}{4m^2} \sum_{i,j=1}^{2m} \binom{2m}{i-j} \left( \sum_a \widehat{\phi}_a^2 \cos^a \frac{(i-j)\pi}{m} \right), \end{aligned}$$

since  $\binom{2m}{i+j} = \binom{2m}{i-j}$ . 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

$$\begin{aligned}
kgk^2 &= \frac{1}{4m^2} \left( 2m \left( \sum_a \widehat{\phi}_a^2 \cos^a 0 \right) + \sum_{t=1}^{2m-1} 2(2m-t)(-1)^t \left( \sum_a \widehat{\phi}_a^2 \cos^a \frac{t\pi}{m} \right) \right) \\
&= \frac{1}{4m^2} \sum_a \widehat{\phi}_a^2 \left( 2m + \sum_{t=1}^{2m-1} 2(2m-t)(-1)^t \cos^a \frac{t\pi}{m} \right) \\
&= \frac{1}{4m^2} \sum_a \widehat{\phi}_a^2 S(a, m), \tag{C.2.1}
\end{aligned}$$

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

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

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

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

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

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

$$= \begin{cases} 2^{1-a} m \left( \sum_{\substack{p=1 \\ p \text{ odd}}}^{ba/mc} \binom{a}{a/2 - pm/2} \right) & a \geq m \\ 0 & a < m \end{cases}$$

To get a sense for the asymptotics as  $a \rightarrow \infty$ , we consider  $a \geq m$  (say  $a \geq 100m$ ). 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. (C.2.1), we get that

$$kgk^2 = \Omega \left( \sum_{\substack{a \leq m \\ a \text{ even}}} \frac{\widehat{\phi}_a^2}{a} \right).$$

□

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

**Corollary C.2.2.** *Consider  $g$  instantiated with  $\phi = \text{ReLU}$ . Then  $kgk = \Omega(1/m)$ .*

*Proof.* The Hermite coefficients of ReLU satisfy  $\widehat{\phi}_a = \Theta(a^{-5/4})$  (Lemma C.3.1). Thus by Lemma C.2.1,

$$kgk^2 = \Omega \left( \sum_{\substack{a \leq 100m \\ a \text{ even}}} a^{-3} \right) = \Omega(1/m^2).$$

□

**Corollary C.2.3.** Consider  $g$  instantiated with  $\phi = \sigma$ , the standard sigmoid. Then  $\|kg\| = e^{-O(\frac{1}{m})}$ .

*Proof.* The Hermite coefficients of  $\sigma$  asymptotically satisfy  $\hat{\phi}_a \leq e^{-C \frac{1}{a}}$  [GGJ+20, §A.2] for some  $C$ . Thus by Lemma C.2.1,

$$\|kg\|^2 = \Omega\left(\sum_{\substack{a=100m \\ a \text{ even}}} e^{-\frac{1}{a}}\right).$$

The result then follows by the following standard integral approximation:

$$\sum_{t=N}^{\infty} \frac{e^{-\frac{1}{t}}}{t} \approx \int_N^{\infty} \frac{e^{-\frac{1}{t}}}{t} dt = 2e^{-\frac{1}{N}}.$$

□

We can now translate these into norm lower bounds on  $f = \psi \circ g$ . For us it suffices to consider  $\psi = \tanh : \mathbb{R} \rightarrow [-1, 1]$ , 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 C.2.4.** Consider  $f$  instantiated with  $\psi = \tanh$  and  $\phi = \text{ReLU}$ . Then  $\|kf\| = \Omega(1/m^6)$ .

*Proof.* Ideally we would like to use the norm bound on  $g$  to obtain an anti-concentration inequality of the form  $\mathbb{P}[|g(x)| > t] \leq \dots$ , 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,  $\text{ReLU}^T(x) = \min(T, \text{ReLU}(x))$  ( $T$  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

$$\mathbb{P}[g(x) \neq g^T(x)] \leq 2m \int_{t=T}^{\infty} \mathbb{P}[\text{ReLU}(t) \neq \text{ReLU}^T(t)] dt \leq 2me^{-T^2/2},$$

since each  $w_i$  is a unit vector.

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

$$\begin{aligned} \|g - g^T\| &= \frac{1}{2m} \sum_{i=1}^{2m} \left( \text{ReLU}_{w_i} - \text{ReLU}_{w_i}^T \right) \\ &= \frac{1}{2m} \sum_{i=1}^{2m} \left\| \text{ReLU}_{w_i} - \text{ReLU}_{w_i}^T \right\| \\ &= \left\| \text{ReLU} - \text{ReLU}^T \right\|_{\mathcal{N}(0,1)} \\ &= \sqrt{e^{-\frac{T^2}{2}} \left( T^2 + 1 + \frac{T}{2\pi} \right)} \end{aligned}$$

where the third equality again uses the fact the  $w_i$  are unit vectors, and the last inequality is Lemma C.2.5. By picking  $T = \Theta(m)$ , this coupled with the fact that  $\|kg\| = \Omega(1/m)$  (Corollary C.2.2) 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:

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

so that

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

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

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

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

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

Thus finally we have

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

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



**Lemma C.2.5** ([GGJ<sup>+</sup>20], Appendix A.1). For  $\text{ReLU}^T(x) = \min(T, \text{ReLU}(x))$ ,

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

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

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

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

$$\begin{aligned} \int_T^1 t^2 p(t) dt &= \int_T^1 t d(p(t)) \\ &= t p(t) \Big|_T^1 + \int_T^1 p(t) dt && \text{(integration by parts)} \\ &= T p(T) + \mathbb{P}_{t \sim N(0,1)}(t > T), \\ \int_T^1 t p(t) dt &= p(t) \Big|_T^1 = p(T), \\ \int_T^1 p(t) dt &= \mathbb{P}_{t \sim N(0,1)}(t > T) = e^{-\frac{T^2}{2}}. \end{aligned}$$

The claim follows by algebra.  $\square$

**Lemma C.2.6.** Consider  $f$  instantiated with  $\psi = \tanh$  and  $\phi = \sigma$ . Then  $\|k\| = e^{-O(\frac{1}{m})}$ .

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

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

In our case, taking  $t = e^{-C/\sqrt{m}}$  for sufficiently large  $C$  and using  $\|kgk\| = e^{-O(\frac{1}{m})}$  (Corollary C.2.3) yields

$$\mathbb{P}[|jg(x)| > e^{-C/\sqrt{m}}] = e^{-O(\frac{1}{m})}.$$

Thus

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

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

### C.3 Approximate degree of ReLUs and sigmoids

Here we give estimates for the  $\delta$ -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  $\|w\|_2 = 1$ . 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  $x \sim N(0, I_n)$  is the same as that of the univariate  $\phi(t)$  as  $t \sim N(0, 1)$ . (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  $\hat{\phi}_a$  are the Hermite coefficients of  $\phi : \mathbb{R} \rightarrow \mathbb{R}$  if  $\phi = \sum_a \hat{\phi}_a H_a$ , where the  $H_a$  are the normalized probabilists' Hermite polynomials. We use  $\tilde{H}_a$  to denote the un-normalized (i.e. monic) Hermite polynomials. (Note that this is somewhat nonstandard notation.)

First we consider ReLUs.

**Lemma C.3.1.**  $\mathbb{E} \text{ReLU}_0 = 1/\sqrt{2\pi}$ ,  $\mathbb{E} \text{ReLU}_1 = 1/2$  and for  $a \geq 2$ ,  $\mathbb{E} \text{ReLU}_a = \frac{1}{2\pi a!} (\tilde{H}_a(0) + a\tilde{H}_{a-2}(0))$ . In particular,  $\mathbb{E} \text{ReLU}_a = 0$  for odd  $a \geq 3$  and  $\mathbb{E} \text{ReLU}_a = \Theta(a^{-5/4})$  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  $a \geq 2$ ,

$$\begin{aligned}
\mathbb{E}\text{ReLU}_a &= \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} \text{ReLU}(x) H_a(x) e^{-\frac{x^2}{2}} dx \\
&= \frac{1}{\sqrt{2\pi} a!} \int_0^{\infty} x \tilde{H}_a(x) e^{-\frac{x^2}{2}} dx \\
&= \frac{1}{\sqrt{2\pi} a!} \int_0^{\infty} (\tilde{H}_{a+1}(x) + a\tilde{H}_{a-1}(x)) e^{-\frac{x^2}{2}} dx \\
&= \frac{1}{\sqrt{2\pi} a!} (\tilde{H}_a(0) + a\tilde{H}_{a-2}(0)).
\end{aligned}$$

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

$$\begin{aligned}
\mathbb{E}\text{ReLU}_a &= \frac{1}{\sqrt{2\pi} (2b)!} (\tilde{H}_{2b}(0) + 2b\tilde{H}_{2b-2}(0)) \\
&= \frac{1}{\sqrt{2\pi} (2b)!} \left( \binom{2b}{b} \frac{(2b)!}{b! 2^b} + 2b \binom{2b-2}{b-1} \frac{(2b-2)!}{(b-1)! 2^{b-1}} \right) \\
&= \frac{\binom{2b-1}{b} \sqrt{(2b)!}}{2\pi b! 2^b} \left( 1 + \frac{2b}{2b-1} \right) \\
&= \frac{\binom{2b-1}{b} \sqrt{(2b)!}}{2\pi (2b-1) b! 2^b} \\
&\leq \frac{\binom{2b-1}{b}}{2\pi (2b-1) (2b)^{1/4}} \\
&\leq \frac{\binom{2b-1}{b}}{b^{5/4}}
\end{aligned}$$

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

**Corollary C.3.2.** *The  $\delta$ -approximate degree of ReLU under  $N(0, 1)$  is  $O(\delta^{-4/3})$ .*

*Proof.* Let  $p$  denote the the Hermite expansion of ReLU truncated at degree  $d$ . By

the fact that  $\sum_{a \geq j} \text{ReLU}_a^j = \Theta(a^{-5/4})$  for even  $a$  (and 0 for odd  $a$ ), we see that

$$\begin{aligned} \sum_{a \geq d} \text{ReLU}_a^2 &= \sum_{a > d} \text{ReLU}_a^2 \\ &= \sum_{\substack{a > d \\ a \text{ even}}} \Theta(a^{-5/2}) \\ &= \Theta(d^{-3/2}). \end{aligned}$$

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

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

**Lemma C.3.3.** *For all sufficiently large  $a$ ,  $\widehat{\sigma}_a = e^{-\binom{p}{a}}$ .*

*Proof.* Upper bounds on the Hermite coefficients of sigmoidal functions 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  $\tanh^\ell(x) = 1 - \tanh^2(x)$ . The calculation is very similar for  $\sigma$  (in fact,  $\sigma$  is just an affine shift of  $\tanh$ ).  $\square$

**Corollary C.3.4.** *The  $\delta$ -approximate degree of  $\sigma$  under  $N(0, 1)$  is  $\tilde{O}(\log^2 1/\delta)$ .*

*Proof.* Let  $p$  denote the Hermite expansion of  $\sigma$  truncated at degree  $d$ . Observe that

$$\begin{aligned} \sum_{a \geq d} p_a^2 &= \sum_{a > d} \widehat{\sigma}_a^2 \\ &= \sum_{a > d} e^{-\binom{p}{a}} \\ &= \Theta\left(\frac{1}{de^{-\binom{p}{d}}}\right), \end{aligned}$$

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

## C.4 Frank–Wolfe convergence guarantee

Here we provide a self-contained proof of Theorem 4.2.5, 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

$$hs, \quad \langle p(z_t), s \rangle \leq \alpha \max_{s' \in Z^0} \langle p(z_t), s' \rangle - \frac{1}{2} \delta \gamma_t C_p \quad (\text{C.4.1})$$

for some constant  $\alpha \geq 1$ . We closely follow [Jag13, Appendix A], noting the differences in our slightly more general setup (the standard setup has  $Z^0 = Z$ , and  $\alpha = 1$ ).

**Theorem C.4.1.** *Let  $Z^0 \subseteq Z$  be convex sets, and let  $p : Z \rightarrow \mathbb{R}$  be a  $\beta$ -smoothly convex function. Let  $C_p = \beta \text{diam}(Z)^2$ . Suppose that  $z \in Z^0$  achieves  $\min_{z' \in Z^0} p(z')$ . For every  $t$ , the iterates of Algorithm 1 (modified to work with Eq. (C.4.1)) satisfy*

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

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

$$q(z) = \max_{s \in Z^0} \langle z, s \rangle - p(z).$$

Notice that  $q$  takes in any  $z \in Z$  but maximizes only over  $s \in Z^0$ . By convexity of  $p$  over  $Z$ , we know that for all  $z \in Z, s \in Z^0, p(z) + \langle z, s \rangle - p(s) \leq q(z)$ , 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  $t^{\text{th}}$  step be  $z_{t+1} = z_t + \gamma(s - z_t)$ , where  $z_t, z_{t+1}, s \in Z, \gamma \in [0, 1]$  is arbitrary, and  $s$  satisfies*

$$hs, \quad \langle p(z_t), s \rangle \leq \alpha \max_{s' \in Z^0} \langle p(z_t), s' \rangle - \frac{1}{2} \delta \gamma C_p.$$

*Then we have*

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

To see this, first note that because  $p$  is  $\beta$ -smoothly convex,

$$p(z_{t+1}) = p(z_t + \gamma(s - z_t)) \\ p(z_t) + \gamma \langle s - z_t, \nabla p(z_t) \rangle + \frac{\gamma^2}{2} C_p.$$

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

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

The claim now follows.

As a consequence of the claim, we can say

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

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

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

This proves the theorem. □

## C.5 Relationship between Boolean 0-1 loss and real-valued correlation loss

Let  $D$  be a distribution on  $\mathbb{R}^n \times \mathbb{R}$ . Our lower bound applies against agnostic learners that satisfy Assumption 4.3.1, with a real-valued correlation guarantee, i.e. learners that learn a class  $H$  by outputting  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  such that

$$\mathbb{E}_{(x,y) \sim D} [f(x)y] \geq \max_{g \in H} \mathbb{E}_{(x,y) \sim D} [g(x)y] - \epsilon. \tag{C.5.1}$$

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

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

or equivalently

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

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

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  $\mathbb{R}^n \times \{-1, 1\}$  (with Boolean labels), not distributions  $D$  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  $\phi : \mathbb{R}^n \times \{-1, 1\} \rightarrow \mathbb{R}$ , into queries to  $D$ . We claim that both of these difficulties can be gotten around. We will show that if  $D$  has bounded labels, say in  $[-C, C]$ , we can construct a distribution  $P$  on  $\mathbb{R}^n \times \{-1, 1\}$  and simulate  $L$  on  $P$  to obtain a correlation loss guarantee wrt  $D$ .

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

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

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

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

so that using  $L$  to learn up to 0-1 error  $\epsilon$  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  $D$ . For any query  $\phi : \mathbb{R}^n \rightarrow \mathbb{R}$ , observe that (since the marginal of  $P$  on  $\mathbb{R}^n$  is also  $D$ )

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

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

In our reduction (Theorem 4.4.1), we end up using the base learner on labeled distributions  $D$  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  $\psi : \mathbb{R} \rightarrow [-1, 1]$ . Recall that in solving the Frank–Wolfe subproblem we needed to worry about simulating SQ access to this  $D$  using only SQ access to the true  $D_{\psi \circ f}$  (see Eq. (4.4.1) and surrounding discussion). Here we actually have a further layer: we need to simulate SQ access to  $P$  using SQ access to  $D$ , 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. (4.4.1).

## C.6 Relationship between square loss and correlation loss for ReLUs

Let  $D$  be a distribution on  $\mathbb{R}^n \rightarrow \mathbb{R}$ , and assume the labels are bounded in  $[-C, C]$ . Our lower bounds apply to agnostic learners that satisfy Assumption 4.3.1, with a guarantee in terms of correlation, where the output hypothesis  $f$  must satisfy

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

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

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



Here we show that for the class of ReLUs,  $H = H_{\text{ReLU}}$ , an agnostic learner  $L$  with a square loss guarantee can be used to satisfy Assumption 4.3.1. 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  $H_{\text{ReLU}}$  so that the maximum norm of any function is 1:

$$H = H_{\text{ReLU}} = \left\{ \frac{1}{\sqrt{2}} \text{ReLU}(u \cdot x) \mid \|u\|_2 = 1 \right\}.$$

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

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

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

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

$$\|h_{\text{sq}}\| = \min_{\lambda \in [0, 1]} \left\| \lambda \frac{h_{\text{sq}}}{\|h_{\text{sq}}\|} - f_{\text{cmf}} \right\|. \quad (\text{C.6.1})$$

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

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

the observation follows immediately.<sup>1</sup> This projection theorem also tells us that  $h_{\text{sq}} = 0$  iff  $f_{\text{cmf}}$  has no projection onto any  $h \in H$ , i.e.  $\langle h, f_{\text{cmf}} \rangle = 0$  for all  $h \in H$ .<sup>2</sup>

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

$$\begin{aligned} \|h_{\text{sq}} - f_{\text{cmf}}\|^2 &= \|h_{\text{sq}} - \frac{\|h_{\text{sq}}\|}{\|h_{\text{sq}}\|} h_{\text{sq}} + \frac{\|h_{\text{sq}}\|}{\|h_{\text{sq}}\|} h_{\text{sq}} - f_{\text{cmf}}\|^2 \\ \Rightarrow \langle h_{\text{sq}} - f_{\text{cmf}}, h_{\text{sq}} - f_{\text{cmf}} \rangle &= \langle h_{\text{sq}} - f_{\text{cmf}}, h_{\text{sq}} - f_{\text{cmf}} \rangle - \langle h_{\text{sq}} - f_{\text{cmf}}, h_{\text{sq}} - f_{\text{cmf}} \rangle. \end{aligned}$$

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

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

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

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

$$\|h - f_{\text{cmf}}\|^2 \leq \|h_{\text{sq}} - f_{\text{cmf}}\|^2 + \epsilon^\ell = \left\| h_{\text{sq}} - \frac{\|h\|}{\|h_{\text{sq}}\|} h_{\text{sq}} + \frac{\|h\|}{\|h_{\text{sq}}\|} h_{\text{sq}} - f_{\text{cmf}} \right\|^2 + \epsilon^\ell.$$

<sup>1</sup>Note that here we are assuming  $\langle h_{\text{sq}}, f_{\text{cmf}} \rangle \neq 0$  WLOG, since otherwise we would consider  $-h_{\text{sq}}$ .

<sup>2</sup>For another way to see this, for any nonzero  $h \in H$ , expand  $\|k\lambda h - f_{\text{cmf}}\|^2 \leq \|k\lambda h_{\text{sq}} - f_{\text{cmf}}\|^2$  and let  $\lambda \rightarrow 0$ .

Some rearrangement gives

$$\begin{aligned} \left\langle h \frac{h}{khk}, f_{\text{cmf}} \right\rangle &= \left\langle h \frac{h_{\text{sq}}}{kh_{\text{sq}}k}, f_{\text{cmf}} \right\rangle = \frac{\epsilon^\ell}{2khk} \\ &= \left\langle hh_{\text{cor}}, f_{\text{cmf}} \right\rangle = \frac{\epsilon^\ell}{2khk}, \end{aligned} \quad (\text{C.6.2})$$

showing that  $h/khk$  is  $\frac{\epsilon^\ell}{2khk}$ -competitive with  $h_{\text{cor}}$ .

But an issue here is that  $khk$  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 C.6.1.** *Let  $h$  be such that  $kh \langle f_{\text{cmf}}, k^2 \rangle \leq kh_{\text{sq}} \langle f_{\text{cmf}}, k^2 \rangle + \epsilon^\ell$ . Suppose  $khk \geq \eta$ . Then  $\langle hh_{\text{cor}}, f_{\text{cmf}} \rangle \leq \frac{\epsilon^\ell}{\epsilon^\ell + 2C\eta}$ . In particular, the 0 function is  $\frac{\epsilon^\ell}{\epsilon^\ell + 2C\eta}$ -competitive with  $h_{\text{cor}}$ .*

*Proof.* By Cauchy–Schwarz,

$$k0 \langle f_{\text{cmf}}, k^2 \rangle \leq kh \langle f_{\text{cmf}}, k^2 \rangle = 2\langle hh, f_{\text{cmf}} \rangle \leq kf_{\text{cmf}}k^2 \leq 2khk kf_{\text{cmf}}k \leq 2C\eta,$$

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

Thus

$$k0 \langle f_{\text{cmf}}, k^2 \rangle \leq kh \langle f_{\text{cmf}}, k^2 \rangle + 2C\eta \leq kh_{\text{sq}} \langle f_{\text{cmf}}, k^2 \rangle + \epsilon^\ell + 2C\eta.$$

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

$$kh_{\text{sq}} \langle f_{\text{cmf}}, k^2 \rangle \leq k0 \langle f_{\text{cmf}}, k^2 \rangle,$$

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

$$kh_{\text{sq}} \langle f_{\text{cmf}}, k^2 \rangle \leq k0 \langle f_{\text{cmf}}, k^2 \rangle \leq kh_{\text{sq}} \langle f_{\text{cmf}}, k^2 \rangle + \epsilon^\ell + 2C\eta. \quad (\text{C.6.3})$$

This lets us conclude that  $kh_{\text{sq}}k$  must be small:

$$kh_{\text{sq}}k^2 = kf_{\text{cmf}}k^2 \leq kh_{\text{sq}} \langle f_{\text{cmf}}, k^2 \rangle + 2\langle hh_{\text{sq}}, f_{\text{cmf}} \rangle \leq \epsilon^\ell + 2C\eta,$$

where we use Eq. (C.6.3) and the fact that by can rewrite Eq. (C.6.1) as  $\langle h, h_{\text{sq}} \rangle = \frac{\langle h, h_{\text{sq}} \rangle}{\|h_{\text{sq}}\|} \|h_{\text{sq}}\|$  or  $\langle h, h_{\text{sq}} \rangle = \langle f_{\text{cmf}}, h_{\text{sq}} \rangle = 0$ . But now since  $\frac{\langle h, h_{\text{sq}} \rangle}{\|h_{\text{sq}}\|} \leq \sqrt{\epsilon^\ell + 2C\eta} < 1$  ( $\epsilon^\ell$  and  $\eta$  will be picked sufficiently small), Eq. (C.6.1) boils down to saying that

$$\langle h, h_{\text{cor}} \rangle = \frac{\langle h, h_{\text{sq}} \rangle}{\|h_{\text{sq}}\|}, \langle f_{\text{cmf}}, h_{\text{cor}} \rangle = \frac{\langle f_{\text{cmf}}, h_{\text{sq}} \rangle}{\|h_{\text{sq}}\|} \sqrt{\epsilon^\ell + 2C\eta}.$$

□

We can now put everything together.

**Theorem C.6.2.** *Suppose we have an agnostic learner  $L$  for  $H_{\text{ReLU}}$  under  $D$  with a square loss guarantee. Then  $L$  can be used to yield a correlation guarantee, i.e. to satisfy Assumption 4.3.1.*

*Proof.* Run  $L$  with  $\epsilon^\ell = \Theta(\epsilon^3)$  to get  $h$  such that  $\langle h, h_{\text{sq}} \rangle = \langle f_{\text{cmf}}, h_{\text{sq}} \rangle + \epsilon^\ell$ . By Lemma C.6.1, if  $\|h\| = \eta = \Theta(\epsilon^2)$ , then  $h$  is  $\epsilon$ -competitive with  $h_{\text{cor}}$ . So we may assume that  $\|h\| = \Theta(\epsilon^2)$ . But then by Eq. (C.6.2), since now  $\frac{\langle h, h_{\text{sq}} \rangle}{\|h\|} \geq \epsilon$ , we get that  $h / \|h\|$  is  $\epsilon$ -competitive with  $h_{\text{cor}}$ . □

# Appendix D: A Moment-Matching Approach to Testable Learning and a New Characterization of Rademacher Complexity

## D.1 Proof of strong duality in Theorem 6.3.2

We will use the following statement of conic duality, specialized to the setting of moment problems.

**Theorem D.1.1** ([Sha01, Section 3]). *Let  $\Omega = \mathbb{R}^d$ , endowed with the standard Borel sigma algebra, and let  $\mathcal{C}$  be the set of all nonnegative Borel measures on  $\Omega$ . Pair the space of signed measures on  $\Omega$  and functions mapping  $\Omega$  to  $\mathbb{R}$  using the following inner product:  $\langle h, \mu \rangle = \int g d\mu$ . Let  $\phi, \psi_1, \dots, \psi_p : \Omega \rightarrow \mathbb{R}$  be functions, let  $b \in \mathbb{R}^p$ , let  $A : \mu \mapsto (\langle \psi_1, \mu \rangle, \dots, \langle \psi_p, \mu \rangle)$ , and let  $K$  be a closed convex cone in  $\mathbb{R}^p$ .*

Define the following primal problem ([Sha01, Eq 3.2]):

$$\sup_{\mu \in \mathcal{C}} \langle \phi, \mu \rangle \quad \text{subject to} \quad A\mu \in b + K. \quad (\text{D.1.1})$$

Let  $K^\circ = \{ \alpha \in \mathbb{R}^p \mid \langle \alpha, b \rangle \leq \langle \alpha, A\mu \rangle \text{ for all } \mu \in \mathcal{C} \}$  be the polar cone of  $K$ . Then the dual is defined as follows ([Sha01, Eq 3.8]):

$$\inf_{\alpha \in K^\circ} \langle b, \alpha \rangle \quad \text{subject to} \quad \sum_{i=1}^p \alpha_i \psi_i(\omega) \leq \phi(\omega) \quad \forall \omega \in \Omega. \quad (\text{D.1.2})$$

Further, a sufficient condition for strong duality to hold (i.e. for both primal and dual to have the same finite optimum) is that  $b$  lie in the interior of the feasible set, i.e.  $b \in \text{int}\{ \mu \in \mathcal{C} : A\mu \in K \}$  ([Sha01, Eq 3.12]).

Let  $\Omega, \mathcal{C}$  and the dual pairing  $\langle \cdot, \cdot \rangle$  be as above. Note that  $\mathcal{C}$  can also be viewed as the convex cone generated by all Dirac measures on  $\Omega$ , and also that when  $\mu$  is a probability measure (i.e., nonnegative and with total measure 1),  $\langle h, \mu \rangle = \mathbb{E}_\mu[h]$ .

Our goal now is to obtain strong duality between Eqs. (6.3.1) and (6.3.3) as a consequence of Theorem D.1.1. Let  $r = j \mid (k, d) j = 1$ , and for convenience write  $l(k, d) = f \mid I_0, I_1, \dots, I_r g$ , where  $I_0 = (0, \dots, 0)$ . Define the functions  $\psi_1, \dots, \psi_r : \Omega \rightarrow \mathbb{R}$  to be the nontrivial monomials corresponding to  $l(k, d)$ , i.e.,  $\psi_j(x) = x_{I_j}$ , and define  $\psi_{r+j} = \psi_r$  for all  $1 \leq j \leq r$ . Let  $K \subset \mathbb{R}^{2r+1}$  be the following convex cone:  $K = f \mid 0 g \subset \mathbb{R}^{2r}$ , where  $\mathbb{R} = (-1, 0]$ . Let  $A$  be a linear operator on  $\mathcal{C}$  given by

$$A\mu = (h \mid 1, \mu i, h \psi_1, \mu i, \dots, h \psi_r, \mu i, h \mid \psi_1, \mu i, \dots, h \mid \psi_r, \mu i) \quad (\text{D.1.3})$$

and let  $b \in \mathbb{R}^{2r+1}$  be given by

$$b = (1, \sigma_{I_1} + \Delta_{I_1}, \dots, \sigma_{I_r} + \Delta_{I_r}, \sigma_{I_1} + \Delta_{I_1}, \dots, \sigma_{I_r} + \Delta_{I_r}). \quad (\text{D.1.4})$$

We claim that our original primal LP Eq. (6.3.1) corresponds to the following conic linear program, which has the form of Eq. (D.1.1), with  $p = 2r + 1$  and  $\phi = f$ :

$$\sup_{\mu \in \mathcal{C}} h \mid f, \mu i \quad \text{subject to} \quad A\mu \leq b \in K \quad (\text{D.1.5})$$

Indeed, the first coordinate of  $b$  ensures the  $I = 0$  constraint, namely that  $h \mid 1, \mu i = 1$  and hence  $\mu$  is a valid probability measure (note that the cone  $\mathcal{C}$  already only consists of nonnegative measures), and the other coordinates ensure that  $\sigma_I \leq \Delta_I = \mathbb{E}_\mu[x_I]$  and  $\sigma_I + \Delta_I$  for every other  $I \in l(k, d) \cap f \mid 0 g$ .

The dual of this program may be written in the form of Eq. (D.1.2) as follows. First introduce dual variables  $\alpha_0 \in \mathbb{R}$  (corresponding to the first constraint), and  $(\alpha_1, \dots, \alpha_{2r}) \in \mathbb{R}^{2r}$  (corresponding to the others), and write  $b = (b_0, \dots, b_{2r})$ . The dual is

$$\inf_{\alpha \in K} \alpha_0 b_0 + \sum_{j=1}^{2r} \alpha_j b_j \quad \text{subject to} \quad \alpha_0 + \sum_{j=1}^{2r} \alpha_j \psi_j \leq f \text{ over } \Omega. \quad (\text{D.1.6})$$

Here  $K$  is the polar cone of  $K$ , and is easily seen to be  $K = \mathbb{R} \times \mathbb{R}^{2r}$ . This means  $K = \mathbb{R} \times \mathbb{R}_+^{2r}$ , i.e.  $\alpha_0 \in \mathbb{R}$  and  $(\alpha_1, \dots, \alpha_{2r}) \in \mathbb{R}_+^{2r}$ . The dual objective may be

simplified as follows:

$$\alpha_0 b_0 + \sum_{j=1}^{2r} \alpha_j b_j = \alpha_0 + \sum_{j=1}^r (\alpha_j (\sigma_{I_j} + \Delta_{I_j}) + \alpha_{r+j} (\sigma_{I_j} + \Delta_{I_j})) \quad (\text{D.1.7})$$

$$= \alpha_0 + \sum_{j=1}^r (\alpha_j + \alpha_{r+j}) \sigma_{I_j} + \sum_{j=1}^r (\alpha_j + \alpha_{r+j}) \Delta_{I_j}. \quad (\text{D.1.8})$$

The constraint simplifies to

$$\alpha_0 + \sum_{j=1}^r (\alpha_j + \alpha_{r+j}) \psi_j \leq f.$$

To simplify this further, if we let  $\beta_j = \alpha_j + \alpha_{r+j}$  for every  $1 \leq j \leq r$ , then it is not hard to see that the objective is minimized when each  $\alpha_j + \alpha_{r+j} = j\beta_j$  (in particular, when  $\alpha_j = \max\{j\beta_j, 0\}$  and  $\alpha_{r+j} = \max\{j\beta_j - \alpha_j, 0\}$ ). Thus if we also let  $\beta_0 = \alpha_0$ , then the dual objective becomes  $\beta_0 + \sum_{j=1}^r \beta_j \sigma_{I_j} + \sum_{j=1}^r j\beta_j \Delta_{I_j}$ , and the constraint becomes  $\beta_0 + \sum_{j=1}^r \beta_j \psi_j \leq f$ . Recalling that  $\sigma_{I_0} = 1$  and  $\Delta_{I_0} = \Delta_0 = 0$ , this is precisely the dual we originally claimed, Eq. (6.3.3).

Now, by Theorem D.1.1, a sufficient condition for strong duality is that  $b$  lie in the interior of the feasible set, i.e.  $b \in \text{int} \{ \int \tilde{b} \, d\mu \in C : A\mu \leq \tilde{b} \in K \}$ . This means that for any sufficiently small perturbation  $\tilde{b}$  of  $b$ , there must exist a measure  $\mu \in C$  such that  $A\mu \leq \tilde{b} \in K$ , i.e. with  $\tilde{b}$  as its approximate vector of moments up to order  $k$ . We argue this slightly informally as follows. Let  $\mu$  denote  $D$  from the statement of Theorem 6.3.2. Suppose

$$\tilde{b} = b + \eta = (b_0 + \eta_0, b_1 + \eta_1, \dots, b_{r+1} + \eta_{r+1}, \dots) \quad (\text{D.1.9})$$

$$= (1 + \eta_0, \sigma_{I_1} + \Delta_{I_1} + \eta_1, \dots, \sigma_{I_1} + \Delta_{I_1} + \eta_{r+1}, \dots), \quad (\text{D.1.10})$$

where  $\eta_0, \dots, \eta_{2r} \in \mathbb{R}^{2r+1}$  are to be thought of as small. The condition that  $A\mu \leq \tilde{b} \in K$  is the same as saying that  $\mu$  satisfies the following:

$$h\mu, 1i = 1 + \eta_0 \quad (\text{D.1.11})$$

$$\sigma_{I_j} + \Delta_{I_j} + \eta_{r+j} \leq h\mu, \psi_ji \leq \sigma_{I_j} + \Delta_{I_j} + \eta_j \quad \forall j = 1, \dots, r. \quad (\text{D.1.12})$$

For sufficiently small  $\eta$ , we claim that a small perturbation of  $\mu$  will continue to satisfy these conditions. First, note that because  $\int \mu, 1 \neq 1$ ,  $\mu$  is no longer formally a probability measure. But for sufficiently small  $\eta_0$ , by adding or removing some mass to  $\mu$  arbitrarily close to the origin, we can increase or decrease its total mass while keeping all its moments nearly unchanged (because the  $\psi_j$  are continuous and  $\psi_j(0) = 0$  for all  $j \neq 0$ , and the new mass is essentially all at 0). Take  $\mu$  to be such a perturbation of  $\mu$ , satisfying  $\int \mu, 1 = 1 + \eta_0$ . We have just argued that for every  $j \neq 0$ ,  $\int \mu, \psi_j$  differs from  $\int \mu, \psi_j = \sigma_{I_j}$  by an arbitrarily small amount. Thus if  $\eta_1, \dots, \eta_{2r}$  are sufficiently small (it suffices to have each  $\eta_j < \Delta_{I_j}/2$ ), then the approximate moment matching conditions will still be satisfied by  $\mu$ , because there is still a slack of at least  $\Delta_{I_j}/2 > 0$  in the constraint arising from  $I_j$ . This establishes that  $b$  is indeed in the interior of the feasible set, and hence that strong duality holds between Eq. (6.3.1) and Eq. (6.3.3).



## Bibliography

- [AAK21] Naman Agarwal, Pranjal Awasthi, and Satyen Kale. A deep conditioning treatment of neural networks. In *Algorithmic Learning Theory*, pages 249–305. PMLR, 2021.
- [ADHV19] Alexandr Andoni, Rishabh Dudeja, Daniel Hsu, and Kiran Vodrahalli. Attribute-efficient learning of monomials over highly-correlated variables. In *Thirtieth International Conference on Algorithmic Learning Theory*, 2019.
- [AGM03] Noga Alon, Oded Goldreich, and Yishay Mansour. Almost k-wise independence versus k-wise independence. *Information Processing Letters*, 88(3):107–110, 2003.
- [AK95] Dana Angluin and Michael Kharitonov. When won’t membership queries help? *Journal of Computer and System Sciences*, 50(2):336–355, 1995.
- [AKPW13] Joël Alwen, Stephan Krenn, Krzysztof Pietrzak, and Daniel Wichs. Learning with rounding, revisited. In *Annual Cryptology Conference*, pages 57–74. Springer, 2013.
- [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.
- [AS20] Emmanuel Abbe and Colin Sandon. Poly-time universality and limitations of deep learning. *arXiv preprint arXiv:2001.02992*, 2020.
- [ATV21] Pranjal Awasthi, Alex Tang, and Aravindan Vijayaraghavan. Efficient algorithms for learning depth-2 neural networks with general relu activations. *Advances in Neural Information Processing Systems*, 34, 2021.

- [Bac17] Francis Bach. Breaking the curse of dimensionality with convex neural networks. *The Journal of Machine Learning Research*, 18(1):629–681, 2017.
- [Bar14] Peter L Bartlett. UC Berkeley CS281B/Stat241B: Statistical Learning Theory, Lecture 7, 2014. URL: <https://www.stat.berkeley.edu/~bartlett/courses/2014spring-cs281bstat241b/lectures/07-notes.pdf>. Last visited on 2022/10/29.
- [Baz09] Louay MJ Bazzi. Polylogarithmic independence can fool dnf formulas. *SIAM Journal on Computing*, 38(6):2220–2272, 2009.
- [BB20] Matthew Brennan and Guy Bresler. Reducibility and statistical-computational gaps from secret leakage. In *Conference on Learning Theory*, pages 648–847. PMLR, 2020.
- [BBH<sup>+</sup>20] Matthew Brennan, Guy Bresler, Samuel B Hopkins, Jerry Li, and Tselil Schramm. Statistical query algorithms and low-degree tests are almost equivalent. *arXiv preprint arXiv:2009.06107*, 2020.
- [BBL02] Peter L Bartlett, Stéphane Boucheron, and Gábor Lugosi. Model selection and error estimation. *Machine Learning*, 48(1):85–113, 2002.
- [BBL03] Olivier Bousquet, Stéphane Boucheron, and Gábor Lugosi. Introduction to statistical learning theory. In *Summer school on machine learning*, pages 169–207. Springer, 2003.
- [BBM05] Peter L Bartlett, Olivier Bousquet, and Shahar Mendelson. Local rademacher complexities. *The Annals of Statistics*, 33(4):1497–1537, 2005.
- [BDES02] Shai Ben-David, Nadav Eiron, and Hans Ulrich Simon. Limitations of learning via embeddings in euclidean half spaces. *Journal of Machine Learning Research*, 3(Nov):441–461, 2002.

- [BDH<sup>+</sup>20] Ainesh Bakshi, Ilias Diakonikolas, Samuel B. Hopkins, Daniel Kane, Sushrut Karmalkar, and Pravesh K. Kothari. Outlier-robust clustering of gaussians and other non-spherical mixtures. In Sandy Irani, editor, *61st IEEE Annual Symposium on Foundations of Computer Science, FOCS 2020, Durham, NC, USA, November 16-19, 2020*, pages 149–159. IEEE, 2020.
- [BdW01] Harry Buhrman and Ronald de Wolf. Communication complexity lower bounds by polynomials. In *Proceedings 16th Annual IEEE Conference on Computational Complexity*, pages 120–130. IEEE, 2001.
- [BDW02] Harry Buhrman and Ronald De Wolf. Complexity measures and decision tree complexity: a survey. *Theoretical Computer Science*, 288(1):21–43, 2002.
- [BE85] D. Bakry and Michel Émery. Diffusions hypercontractives. In *Seminaire de probabilites, XIX, 1983/84*, volume 1123 of *Lecture Notes in Math.*, pages 177–206. Springer, Berlin, 1985.
- [Bel21] Mikhail Belkin. Fit without fear: remarkable mathematical phenomena of deep learning through the prism of interpolation. *Acta Numerica*, 30:203–248, 2021.
- [BF02] Nader H Bshouty and Vitaly Feldman. On using extended statistical queries to avoid membership queries. *Journal of Machine Learning Research*, 2(Feb):359–395, 2002.
- [BFJ<sup>+</sup>94] Avrim Blum, Merrick Furst, Jeffrey Jackson, Michael Kearns, Yishay Mansour, and Steven Rudich. Weakly learning dnf and characterizing statistical query learning using fourier analysis. In *Proceedings of the twenty-sixth annual ACM symposium on Theory of computing*, pages 253–262, 1994.

- [BG17a] Alon Brutzkus and Amir Globerson. Globally optimal gradient descent for a convnet with gaussian inputs. *CoRR*, abs/1702.07966, 2017.
- [BG17b] Alon Brutzkus and Amir Globerson. Globally optimal gradient descent for a convnet with gaussian inputs. In *Proceedings of the 34th International Conference on Machine Learning-Volume 70*, pages 605–614, 2017.
- [BGM<sup>+</sup>16] Andrej Bogdanov, Siyao Guo, Daniel Masny, Silas Richelson, and Alon Rosen. On the hardness of learning with rounding over small modulus. In *Theory of Cryptography Conference*, pages 209–224. Springer, 2016.
- [BGML<sup>+</sup>18] Sauvik Bhattacharya, Oscar Garcia-Morchon, Thijs Laarhoven, Ronald Rietman, Markku-Juhani O Saarinen, Ludo Tolhuizen, and Zhenfei Zhang. Round5: Compact and fast post-quantum public-key encryption. *IACR Cryptol. ePrint Arch.*, 2018:725, 2018.
- [BHKL15] Alina Beygelzimer, Elad Hazan, Satyen Kale, and Haipeng Luo. Online gradient boosting. In *Advances in neural information processing systems*, pages 2458–2466, 2015.
- [BI91] Gyora M Benedek and Alon Itai. Learnability with respect to fixed distributions. *Theoretical Computer Science*, 86(2):377–389, 1991.
- [BIP<sup>+</sup>18] Dan Boneh, Yuval Ishai, Alain Passelègue, Amit Sahai, and David J Wu. Exploring crypto dark matter. In *Theory of Cryptography Conference*, pages 699–729. Springer, 2018.
- [BJW19] Ainesh Bakshi, Rajesh Jayaram, and David P Woodruff. Learning two layer rectified neural networks in polynomial time. In *Conference on Learning Theory*, pages 195–268. PMLR, 2019.
- [BK20] Ainesh Bakshi and Pravesh Kothari. Outlier-robust clustering of non-spherical mixtures. *CoRR*, abs/2005.02970, 2020.

- [BK21] Ainesh Bakshi and Pravesh K. Kothari. List-decodable subspace recovery: Dimension independent error in polynomial time. In Dániel Marx, editor, *Proceedings of the 2021 ACM-SIAM Symposium on Discrete Algorithms, SODA 2021, Virtual Conference, January 10 - 13, 2021*, pages 1279–1297. SIAM, 2021.
- [BL97] Avrim L Blum and Pat Langley. Selection of relevant features and examples in machine learning. *Artificial intelligence*, 97(1-2):245–271, 1997.
- [BLLT20] Peter L Bartlett, Philip M Long, Gábor Lugosi, and Alexander Tsigler. Benign overfitting in linear regression. *Proceedings of the National Academy of Sciences*, 117(48):30063–30070, 2020.
- [BLMT23] Guy Blanc, Jane Lange, Ali Malik, and Li-Yang Tan. Lifting uniform learners via distributional decomposition. In *Proceedings of the 55th Annual ACM Symposium on Theory of Computing*, pages 1755–1767, 2023.
- [Blu90] Avrim Blum. Learning boolean functions in an infinite attribute space. In *Proceedings of the twenty-second annual ACM symposium on Theory of computing*, pages 64–72, 1990.
- [BM97] Lucien Birgé and Pascal Massart. From model selection to adaptive estimation. In *Festschrift for lucien le cam*, pages 55–87. Springer, 1997.
- [BM02] Peter L Bartlett and Shahar Mendelson. Rademacher and gaussian complexities: Risk bounds and structural results. *Journal of Machine Learning Research*, 3(Nov):463–482, 2002.
- [BMR21] Peter L Bartlett, Andrea Montanari, and Alexander Rakhlin. Deep learning: a statistical viewpoint. *Acta numerica*, 30:87–201, 2021.

- [Bog21] Andrej Bogdanov. Personal communication, 2021.
- [Boy84] John P Boyd. Asymptotic coefficients of hermite function series. *Journal of Computational Physics*, 54(3):382–410, 1984.
- [BP14] Abhishek Banerjee and Chris Peikert. New and improved key-homomorphic pseudorandom functions. In *Annual Cryptology Conference*, pages 353–370. Springer, 2014.
- [BPR12] Abhishek Banerjee, Chris Peikert, and Alon Rosen. Pseudorandom functions and lattices. In *Annual International Conference on the Theory and Applications of Cryptographic Techniques*, pages 719–737. Springer, 2012.
- [BR89] Avrim Blum and Ronald L Rivest. Training a 3-node neural network is NP-complete. In *Advances in neural information processing systems*, pages 494–501, 1989.
- [BR17] Andrej Bogdanov and Alon Rosen. Pseudorandom functions: Three decades later. In *Tutorials on the Foundations of Cryptography*, pages 79–158. Springer, 2017.
- [Bra10] Mark Braverman. Polylogarithmic independence fools  $AC^0$  circuits. *Journal of the ACM (JACM)*, 57(5):1–10, 2010.
- [BRST21] Joan Bruna, Oded Regev, Min Jae Song, and Yi Tang. Continuous lwe. In *Proceedings of the 53rd Annual ACM SIGACT Symposium on Theory of Computing*, pages 694–707, 2021.
- [BT96] Nader H Bshouty and Christino Tamon. On the fourier spectrum of monotone functions. *Journal of the ACM (JACM)*, 43(4):747–770, 1996.

- [BT19] Mark Bun and Justin Thaler. A Nearly Optimal Lower Bound on the Approximate Degree of  $AC^0$ . *SIAM Journal on Computing*, 49(4):FOCS17–59–FOCS17–96, 2019.
- [BW08] Peter L Bartlett and Marten H Wegkamp. Classification with a reject option using a hinge loss. *Journal of Machine Learning Research*, 9(8), 2008.
- [CDM16] Corinna Cortes, Giulia DeSalvo, and Mehryar Mohri. Learning with rejection. In *Algorithmic Learning Theory: 27th International Conference, ALT 2016, Bari, Italy, October 19-21, 2016, Proceedings 27*, pages 67–82. Springer, 2016.
- [CGMK22] Sitan Chen, Aravind Gollakota, Raghu Meka, and Adam Klivans. Hardness of noise-free learning for two-hidden-layer neural networks. *Advances in Neural Information Processing Systems*, 2022.
- [CGV15] Aloni Cohen, Shafi Goldwasser, and Vinod Vaikuntanathan. Aggregate pseudorandom functions and connections to learning. In *Theory of Cryptography Conference*, pages 61–89. Springer, 2015.
- [CKLS18] Jung Hee Cheon, Duhyeong Kim, Joohee Lee, and Yongsoo Song. Lizard: Cut off the tail! a practical post-quantum public-key encryption from lwe and lwr. In *International Conference on Security and Cryptography for Networks*, pages 160–177. Springer, 2018.
- [CKM20] Sitan Chen, Adam R Klivans, and Raghu Meka. Learning deep relu networks is fixed-parameter tractable. *arXiv preprint arXiv:2009.13512*, 2020.
- [CKM21] Sitan Chen, Adam Klivans, and Raghu Meka. Efficiently learning one hidden layer relu networks from queries. In *Advances in Neural Information Processing Systems*, 2021.

- [CST<sup>+</sup>00] Nello Cristianini, John Shawe-Taylor, et al. *An introduction to support vector machines and other kernel-based learning methods*. Cambridge university press, 2000.
- [Dan16a] Amit Daniely. Complexity theoretic limitations on learning halfspaces. In *Proceedings of the forty-eighth annual ACM symposium on Theory of Computing*, pages 105–117, 2016.
- [Dan16b] Amit Daniely. Complexity theoretic limitations on learning halfspaces. In *Proceedings of the forty-eighth annual ACM symposium on Theory of Computing*, pages 105–117, 2016.
- [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.
- [DG21] Amit Daniely and Elad Granot. An exact poly-time membership-queries algorithm for extraction a three-layer relu network. *arXiv preprint arXiv:2105.09673*, 2021.
- [DGJ<sup>+</sup>10] Ilias Diakonikolas, Parikshit Gopalan, Ragesh Jaiswal, Rocco A Servedio, and Emanuele Viola. Bounded independence fools halfspaces. *SIAM Journal on Computing*, 39(8):3441–3462, 2010.
- [DGK<sup>+</sup>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.
- [DGKP20] Abhimanyu Das, Sreenivas Gollapudi, Ravi Kumar, and Rina Panigrahy. On the learnability of random deep networks. In *Proceedings of the Fourteenth Annual ACM-SIAM Symposium on Discrete Algorithms*, pages 398–410. SIAM, 2020.



- [DK20] Ilias Diakonikolas and Daniel M. Kane. Small covers for near-zero sets of polynomials and learning latent variable models. In *2020 IEEE 61st Annual Symposium on Foundations of Computer Science (FOCS)*, pages 184–195, 2020.
- [DK21] Ilias Diakonikolas and Daniel M. Kane. Non-gaussian component analysis via lattice basis reduction, 2021.
- [DKK<sup>+</sup>23] Ilias Diakonikolas, Daniel M Kane, Vasilis Kontonis, Sihan Liu, and Nikos Zarifis. Efficient testable learning of halfspaces with adversarial label noise. *arXiv preprint arXiv:2303.05485*, 2023.
- [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.
- [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.
- [DKPZ21] Ilias Diakonikolas, Daniel M Kane, Thanasis Pittas, and Nikos Zarifis. The optimality of polynomial regression for agnostic learning under gaussian marginals in the sq model. In *Conference on Learning Theory*, pages 1552–1584. PMLR, 2021.
- [DKR23] Ilias Diakonikolas, Daniel Kane, and Lisheng Ren. Near-optimal cryptographic hardness of agnostically learning halfspaces and relu regression under gaussian marginals. In *International Conference on Machine Learning*, pages 7922–7938. PMLR, 2023.
- [DKRV18] Jan-Pieter D’Anvers, Angshuman Karmakar, Sujoy Sinha Roy, and Frederik Vercauteren. Saber: Module-lwr based key exchange, cpa-secure

- encryption and cca-secure kem. In *International Conference on Cryptology in Africa*, pages 282–305. Springer, 2018.
- [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.
- [DKZ20a] Ilias Diakonikolas, Daniel Kane, and Nikos Zarifis. Near-optimal sq lower bounds for agnostically learning halfspaces and relus under gaussian marginals. *Advances in Neural Information Processing Systems*, 33:13586–13596, 2020.
- [DKZ20b] 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.
- [DLSS14] Amit Daniely, Nati Linial, and Shai Shalev-Shwartz. From average case complexity to improper learning complexity. In *Proceedings of the forty-sixth annual ACM symposium on Theory of computing*, pages 441–448, 2014.
- [DNS22] Anindya De, Shivam Nadimpalli, and Rocco Servedio. Convex influences. *Innovations in Theoretical Computer Science*, 2022.
- [DS16] Amit Daniely and Shai Shalev-Shwartz. Complexity theoretic limitations on learning dnf’s. In *Conference on Learning Theory*, pages 815–830, 2016.
- [DSFT<sup>+</sup>14] Dana Dachman-Soled, Vitaly Feldman, Li-Yang Tan, Andrew Wan, and Karl Wimmer. Approximate resilience, monotonicity, and the complexity of agnostic learning. In *Proceedings of the twenty-sixth annual*

- ACM-SIAM symposium on Discrete algorithms*, pages 498–511. SIAM, 2014.
- [DSS16] Amit Daniely and Shai Shalev-Shwartz. Complexity theoretic limitations on learning dnf’s. In *Conference on Learning Theory*, pages 815–830. PMLR, 2016.
- [DV20a] Amit Daniely and Gal Vardi. Hardness of learning neural networks with natural weights. *arXiv preprint arXiv:2006.03177*, 2020.
- [DV20b] Amit Daniely and Gal Vardi. Hardness of learning neural networks with natural weights. *Advances in Neural Information Processing Systems*, 33, 2020.
- [DV21] Amit Daniely and Gal Vardi. From local pseudorandom generators to hardness of learning. In *Conference on Learning Theory*, pages 1358–1394. PMLR, 2021.
- [Ear19] Mike Earnest. Proving an identity involving the alternating sum of products of binomial coefficients. Mathematics Stack Exchange, 2019. URL: <https://math.stackexchange.com/q/3108805> (version: 2019-02-11).
- [FCG20] Spencer Frei, Yuan Cao, and Quanquan Gu. Agnostic learning of a single neuron with gradient descent. *arXiv preprint arXiv:2005.14426*, 2020.
- [Fel08] Vitaly Feldman. Evolvability from learning algorithms. In *Proceedings of the fortieth annual ACM symposium on Theory of computing*, pages 619–628, 2008.
- [Fel09] Vitaly Feldman. On the power of membership queries in agnostic learning. *The Journal of Machine Learning Research*, 10:163–182, 2009.

- [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.
- [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.
- [Fel14] Vitaly Feldman. Open problem: The statistical query complexity of learning sparse halfspaces. In *Conference on Learning Theory*, pages 1283–1289, 2014.
- [Fel16] Vitaly Feldman. Statistical query learning. In Ming-Yang Kao, editor, *Encyclopedia of Algorithms*, pages 2090–2095. Springer New York, New York, NY, 2016.
- [Fel17] Vitaly Feldman. A general characterization of the statistical query complexity. In *Conference on Learning Theory*, pages 785–830, 2017.
- [FGR<sup>+</sup>17] Vitaly Feldman, Elena Grigorescu, Lev Reyzin, Santosh S Vempala, and Ying Xiao. Statistical algorithms and a lower bound for detecting planted cliques. *Journal of the ACM (JACM)*, 64(2):8, 2017.
- [FGRW12] Vitaly Feldman, Venkatesan Guruswami, Prasad Raghavendra, and Yi Wu. Agnostic learning of monomials by halfspaces is hard. *SIAM Journal on Computing*, 41(6):1558–1590, 2012.
- [Fri01] Jerome H Friedman. Greedy function approximation: a gradient boosting machine. *Annals of statistics*, pages 1189–1232, 2001.
- [Gab81] Yu R Gabovich. Stability of the characterization of the multivariate normal distribution in the skitovich-darmois theorem. *Journal of Soviet Mathematics*, 16(5):1341–1349, 1981.

- [GGJ<sup>+</sup>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.
- [GGK20] Surbhi Goel, Aravind Gollakota, and Adam Klivans. Statistical-query lower bounds via functional gradients. *Advances in Neural Information Processing Systems*, 33:2147–2158, 2020.
- [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.
- [GKK23] Aravind Gollakota, Adam R Klivans, and Pravesh K Kothari. A moment-matching approach to testable learning and a new characterization of rademacher complexity. In *Proceedings of the 55th Annual ACM Symposium on Theory of Computing*, pages 1657–1670, 2023.
- [GKKM20] Shafi Goldwasser, Adam Tauman Kalai, Yael Kalai, and Omar Montasser. Beyond perturbations: Learning guarantees with arbitrary adversarial test examples. *Advances in Neural Information Processing Systems*, 33:15859–15870, 2020.
- [GKKT17a] 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.
- [GKKT17b] Surbhi Goel, Varun Kanade, Adam R. Klivans, and Justin Thaler. Reliably learning the relu in polynomial time. In *COLT*, pages 1004–1042, 2017.

- [GKM18] Surbhi Goel, Adam R. Klivans, and Raghu Meka. Learning one convolutional layer with overlapping patches. In *ICML*, volume 80, pages 1778–1786. PMLR, 2018.
- [GKSV23a] Aravind Gollakota, Adam R Klivans, Konstantinos Stavropoulos, and Arsen Vasilyan. An efficient tester-learner for halfspaces. *arXiv preprint arXiv:2302.14853*, 2023.
- [GKSV23b] Aravind Gollakota, Adam R Klivans, Konstantinos Stavropoulos, and Arsen Vasilyan. Tester-learners for halfspaces: Universal algorithms. *arXiv preprint arXiv:2305.11765*, 2023.
- [GLM18a] Rong Ge, Jason D. Lee, and Tengyu Ma. Learning one-hidden-layer neural networks with landscape design. In *ICLR*. OpenReview.net, 2018.
- [GLM18b] Rong Ge, Jason D Lee, and Tengyu Ma. Learning one-hidden-layer neural networks with landscape design. In *6th International Conference on Learning Representations, ICLR 2018*, 2018.
- [GR09] Venkatesan Guruswami and Prasad Raghavendra. Hardness of learning halfspaces with noise. *SIAM Journal on Computing*, 39(2):742–765, 2009.
- [GRSY21] Shafi Goldwasser, Guy N Rothblum, Jonathan Shafer, and Amir Yehudayoff. Interactive proofs for verifying machine learning. In *12th Innovations in Theoretical Computer Science Conference (ITCS 2021)*. Schloss Dagstuhl-Leibniz-Zentrum für Informatik, 2021.
- [Haz16] Elad Hazan. Introduction to online convex optimization. *Foundations and Trends in Optimization*, 2(3-4):157–325, 2016.

- [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.
- [HMP<sup>+</sup>93] András Hajnal, Wolfgang Maass, Pavel Pudlák, Mario Szegedy, and György Turán. Threshold circuits of bounded depth. *Journal of Computer and System Sciences*, 46(2):129–154, 1993.
- [HMRT22] Trevor Hastie, Andrea Montanari, Saharon Rosset, and Ryan J Tibshirani. Surprises in high-dimensional ridgeless least squares interpolation. *The Annals of Statistics*, 50(2):949–986, 2022.
- [HS07] Lisa Hellerstein and Rocco A Servedio. On pac learning algorithms for rich boolean function classes. *Theoretical Computer Science*, 384(1):66–76, 2007.
- [HS19] Prahladh Harsha and Srikanth Srinivasan. On polynomial approximations to  $AC^0$ . *Random Structures & Algorithms*, 54(2):289–303, 2019.
- [Hua19] Hao Huang. Induced subgraphs of hypercubes and a proof of the sensitivity conjecture. *Annals of Mathematics*, 190(3):949–955, 2019.
- [IK22] Misha Ivkov and Pravesh K. Kothari. List-decodable covariance estimation. In Stefano Leonardi and Anupam Gupta, editors, *STOC '22: 54th Annual ACM SIGACT Symposium on Theory of Computing, Rome, Italy, June 20 - 24, 2022*, pages 1276–1283. ACM, 2022.
- [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.
- [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.

- [JCB<sup>+</sup>20] Matthew Jagielski, Nicholas Carlini, David Berthelot, Alex Kurakin, and Nicolas Papernot. High accuracy and high fidelity extraction of neural networks. In Srdjan Capkun and Franziska Roesner, editors, *29th USENIX Security Symposium, USENIX Security 2020, August 12-14, 2020*, pages 1345–1362. USENIX Association, 2020.
- [JHG18] Arthur Jacot, Clément Hongler, and Franck Gabriel. Neural tangent kernel: Convergence and generalization in neural networks. In Samy Bengio, Hanna M. Wallach, Hugo Larochelle, Kristen Grauman, Nicolò Cesa-Bianchi, and Roman Garnett, editors, *NeurIPS*, pages 8580–8589, 2018.
- [JSA15] Majid Janzamin, Hanie Sedghi, and Anima Anandkumar. Beating the perils of non-convexity: Guaranteed training of neural networks using tensor methods. *arXiv preprint arXiv:1506.08473*, 2015.
- [JWZ20] Rajesh Jayaram, David P. Woodruff, and Qiuyi Zhang. Span recovery for deep neural networks with applications to input obfuscation. In *ICLR*. OpenReview.net, 2020.
- [JZ16] Zhengzhong Jin and Yunlei Zhao. Optimal key consensus in presence of noise. *arXiv preprint arXiv:1611.06150*, 2016.
- [Kan11] Daniel M Kane. The gaussian surface area and noise sensitivity of degree-d polynomial threshold functions. *computational complexity*, 20(2):389–412, 2011.
- [Kea98] Michael Kearns. Efficient noise-tolerant learning from statistical queries. *Journal of the ACM (JACM)*, 45(6):983–1006, 1998.
- [Kha93] Michael Kharitonov. Cryptographic hardness of distribution-specific learning. In *Proceedings of the twenty-fifth annual ACM symposium on Theory of computing*, pages 372–381, 1993.



- [Kha95] Michael Kharitonov. Cryptographic lower bounds for learnability of boolean functions on the uniform distribution. *Journal of Computer and System Sciences*, 50(3):600–610, 1995.
- [KK09] Varun Kanade and Adam Kalai. Potential-based agnostic boosting. In *Advances in neural information processing systems*, pages 880–888, 2009.
- [KK14a] Adam Klivans and Pravesh Kothari. Embedding hard learning problems into gaussian space. In *Approximation, Randomization, and Combinatorial Optimization. Algorithms and Techniques (APPROX/RANDOM)*, 2014.
- [KK14b] Adam R. Klivans and Pravesh Kothari. Embedding hard learning problems into gaussian space. In *APPROX-RANDOM*, volume 28 of *LIPICs*, pages 793–809. Schloss Dagstuhl - Leibniz-Zentrum fuer Informatik, 2014.
- [KK21] Adam Tauman Kalai and Varun Kanade. Efficient learning with arbitrary covariate shift. In *Algorithmic Learning Theory*, pages 850–864. PMLR, 2021.
- [KKK19] Sushrut Karmalkar, Adam R. Klivans, and Pravesh Kothari. List-decodable linear regression. In Hanna M. Wallach, Hugo Larochelle, Alina Beygelzimer, Florence d’Alché-Buc, Emily B. Fox, and Roman Garnett, editors, *Advances in Neural Information Processing Systems 32: Annual Conference on Neural Information Processing Systems 2019, NeurIPS 2019, December 8-14, 2019, Vancouver, BC, Canada*, pages 7423–7432, 2019.
- [KKM12] Adam Tauman Kalai, Varun Kanade, and Yishay Mansour. Reliable agnostic learning. *Journal of Computer and System Sciences*, 78(5):1481–1495, 2012.

- [KKM13] Daniel Kane, Adam Klivans, and Raghu Meka. Learning halfspaces under log-concave densities: Polynomial approximations and moment matching. In *Conference on Learning Theory*, pages 522–545. PMLR, 2013.
- [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.
- [KL01] Matthias Krause and Stefan Lucks. Pseudorandom functions in in  $tc_0$  and cryptographic limitations to proving lower bounds. *computational complexity*, 10(4):297–313, 2001.
- [KL18] Pravesh K Kothari and Roi Livni. Improper learning by refuting. In *9th Innovations in Theoretical Computer Science Conference (ITCS 2018)*. Schloss Dagstuhl-Leibniz-Zentrum fuer Informatik, 2018.
- [KLS95] Ravi Kannan, László Lovász, and Miklós Simonovits. Isoperimetric problems for convex bodies and a localization lemma. *Discrete & Computational Geometry*, 13:541–559, 1995.
- [KM13] Adam Klivans and Raghu Meka. Moment-matching polynomials. *arXiv preprint arXiv:1301.0820*, 2013.
- [KMS20] Pritish Kamath, Omar Montasser, and Nathan Srebro. Approximate is good enough: Probabilistic variants of dimensional and margin complexity. In *Conference on Learning Theory*, 2020.
- [Kol01] Vladimir Koltchinskii. Rademacher penalties and structural risk minimization. *IEEE Transactions on Information Theory*, 47(5):1902–1914, 2001.

- [Kol06] Vladimir Koltchinskii. Local rademacher complexities and oracle inequalities in risk minimization. *The Annals of Statistics*, 34(6):2593–2656, 2006.
- [KOS08] Adam R Klivans, Ryan O’Donnell, and Rocco A Servedio. Learning geometric concepts via gaussian surface area. In *2008 49th Annual IEEE Symposium on Foundations of Computer Science*, pages 541–550. IEEE, 2008.
- [KP00] Vladimir Koltchinskii and Dmitriy Panchenko. Rademacher processes and bounding the risk of function learning. In *High dimensional probability II*, pages 443–457. Springer, 2000.
- [KR96] LB Klebanov and ST Rachev. Proximity of probability measures with common marginals in a finite number of directions. *Lecture Notes-Monograph Series*, pages 162–174, 1996.
- [KR21] Gil Kur and Alexander Rakhlin. On the minimal error of empirical risk minimization. In *Conference on Learning Theory*, pages 2849–2852. PMLR, 2021.
- [KS94a] 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.
- [KS94b] 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.
- [KS04] Adam R Klivans and Rocco A Servedio. Learning DNF in time  $2^{O(n^{1/3})}$ . *Journal of Computer and System Sciences*, 68(2):303–318, 2004.

- [KS07] Adam R Klivans and Alexander A Sherstov. Unconditional lower bounds for learning intersections of halfspaces. *Machine Learning*, 69(2-3):97–114, 2007.
- [KS09] Adam R Klivans and Alexander A Sherstov. Cryptographic hardness for learning intersections of halfspaces. *Journal of Computer and System Sciences*, 75(1):2–12, 2009.
- [KS10] Adam R Klivans and Alexander A Sherstov. Lower bounds for agnostic learning via approximate rank. *Computational Complexity*, 19(4):581–604, 2010.
- [KSP20] Rohan Karthikeyan, Siddharth Sinha, and Vallabh Patil. On the resolution of the sensitivity conjecture. *Bulletin of the American Mathematical Society*, 2020.
- [KSS92] Michael J Kearns, Robert E Schapire, and Linda M Sellie. Toward efficient agnostic learning. In *Proceedings of the 5th annual workshop on Computational learning theory*, pages 341–352, 1992.
- [KSS94] Michael J Kearns, Robert E Schapire, and Linda M Sellie. Toward efficient agnostic learning. *Machine Learning*, 17(2-3):115–141, 1994.
- [KV94a] Michael Kearns and Leslie Valiant. Cryptographic limitations on learning boolean formulae and finite automata. *Journal of the ACM (JACM)*, 41(1):67–95, 1994.
- [KV94b] Michael J Kearns and Umesh Vazirani. *An introduction to computational learning theory*. MIT press, 1994.
- [Led01] Michel Ledoux. *The concentration of measure phenomenon*. Number 89. American Mathematical Soc., 2001.

- [LLL82] Arjen K Lenstra, Hendrik Willem Lenstra, and László Lovász. Factoring polynomials with rational coefficients. *Mathematische annalen*, 261:515–534, 1982.
- [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.
- [LMSS07] Nati Linial, Shahar Mendelson, Gideon Schechtman, and Adi Shraibman. Complexity measures of sign matrices. *Combinatorica*, 27(4):439–463, 2007.
- [LMZ20] Yuanzhi Li, Tengyu Ma, and Hongyang R. Zhang. Learning over-parametrized two-layer neural networks beyond ntk. In *Conference on Learning Theory 2020*, volume 125, pages 2613–2682. PMLR, 2020.
- [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.
- [LY17] Yuanzhi Li and Yang Yuan. Convergence analysis of two-layer neural networks with relu activation. In *Advances in Neural Information Processing Systems 30*, pages 597–607, 2017.
- [MBBF00] Llew Mason, Jonathan Baxter, Peter L Bartlett, and Marcus R Freen. Boosting algorithms as gradient descent. In *Advances in neural information processing systems*, pages 512–518, 2000.
- [MOS03] Elchanan Mossel, Ryan O’Donnell, and Rocco P Servedio. Learning juntas. In *Proceedings of the thirty- fth annual ACM symposium on Theory of computing*, pages 206–212, 2003.
- [MP69] M.L. Minsky and S. Papert. *Perceptrons; an Introduction to Computational Geometry*. MIT Press, 1969.

- [MS22] Eran Malach and Shai Shalev-Shwartz. When hardness of approximation meets hardness of learning. *The Journal of Machine Learning Research*, 23(1):3942–3965, 2022.
- [MSDH19] Smitha Milli, Ludwig Schmidt, Anca D. Dragan, and Moritz Hardt. Model reconstruction from model explanations. In *FAT*, pages 1–9. ACM, 2019.
- [NK19] Vaishnavh Nagarajan and J Zico Kolter. Uniform convergence may be unable to explain generalization in deep learning. *Advances in Neural Information Processing Systems*, 32, 2019.
- [NR97] Moni Naor and Omer Reingold. Number-theoretic constructions of efficient pseudo-random functions. In *Proceedings 38th Annual Symposium on Foundations of Computer Science*, pages 458–467. IEEE, 1997.
- [NS94] Noam Nisan and Mario Szegedy. On the degree of boolean functions as real polynomials. *Computational complexity*, 4(4):301–313, 1994.
- [Pat92] Ramamohan Paturi. On the degree of polynomials that approximate symmetric boolean functions. In *Proceedings of the twenty-fourth annual ACM symposium on Theory of computing*, pages 468–474, 1992.
- [Pei16] Chris Peikert. A decade of lattice cryptography. *Found. Trends Theor. Comput. Sci.*, 10(4):283–424, mar 2016.
- [PMG<sup>+</sup>17] Nicolas Papernot, Patrick D. McDaniel, Ian J. Goodfellow, Somesh Jha, Z. Berkay Celik, and Ananthram Swami. Practical black-box attacks against machine learning. In Ramesh Karri, Ozgur Sinanoglu, Ahmad-Reza Sadeghi, and Xun Yi, editors, *Proceedings of the 2017 ACM on Asia Conference on Computer and Communications Security, AsiaCCS 2017, Abu Dhabi, United Arab Emirates, April 2-6, 2017*, pages 506–519. ACM, 2017.

- [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.
- [PSP17] PSPACEhard. Alternating sum of binomial coefficients identity. Mathematics Stack Exchange, 2017. URL: <https://math.stackexchange.com/q/2183223> (version: 2017-03-12).
- [Raz92] Alexander A Razborov. On small depth threshold circuits. In *Scandinavian Workshop on Algorithm Theory*, pages 42–52. Springer, 1992.
- [Reg09] Oded Regev. On lattices, learning with errors, random linear codes, and cryptography. *Journal of the ACM (JACM)*, 56(6):1–40, 2009.
- [Reg10] Oded Regev. The learning with errors problem. *Invited survey in CCC*, 7(30):11, 2010.
- [Rey20] Lev Reyzin. Statistical queries and statistical algorithms: Foundations and applications. *arXiv preprint arXiv:2004.00557*, 2020.
- [RK20] David Rolnick and Konrad P. Kording. Reverse-engineering deep relu networks. In *ICML*, volume 119 of *Proceedings of Machine Learning Research*, pages 8178–8187. PMLR, 2020.
- [RKSF13] Svetlozar T Rachev, Lev B Klebanov, Stoyan V Stoyanov, and Frank Fabozzi. *The methods of distances in the theory of probability and statistics*, volume 10. Springer, 2013.
- [Rob55] Herbert Robbins. A remark on stirling’s formula. *The American mathematical monthly*, 62(1):26–29, 1955.
- [RR97] Alexander A Razborov and Steven Rudich. Natural proofs. *Journal of Computer and System Sciences*, 55(1):24–35, 1997.

- [RS10] Alexander A Razborov and Alexander A Sherstov. The sign-rank of  $ac^0$ . *SIAM Journal on Computing*, 39(5):1833–1855, 2010.
- [RV22] Ronitt Rubinfeld and Arsen Vasilyan. Testing distributional assumptions of learning algorithms. *Proceedings of the 54th annual ACM Symposium on Theory of Computing*, 2022. To appear.
- [RV23] Ronitt Rubinfeld and Arsen Vasilyan. Testing distributional assumptions of learning algorithms. In *Proceedings of the 55th Annual ACM Symposium on Theory of Computing*, pages 1643–1656, 2023.
- [RY20a] Prasad Raghavendra and Morris Yau. List decodable learning via sum of squares. In Shuchi Chawla, editor, *Proceedings of the 2020 ACM-SIAM Symposium on Discrete Algorithms, SODA 2020, Salt Lake City, UT, USA, January 5-8, 2020*, pages 161–180. SIAM, 2020.
- [RY20b] Prasad Raghavendra and Morris Yau. List decodable subspace recovery, 2020.
- [SB14] Shai Shalev-Shwartz and Shai Ben-David. *Understanding machine learning: From theory to algorithms*. Cambridge university press, 2014.
- [SF12] Robert E. Schapire and Yoav Freund. *Boosting: Foundations and Algorithms*. The MIT Press, 2012.
- [Sha01] Alexander Shapiro. On duality theory of conic linear problems. In *Semi-in nite programming*, pages 135–165. Springer, 2001.
- [Sha18a] Ohad Shamir. Distribution-specific hardness of learning neural networks. *J. Mach. Learn. Res*, 19:32:1–32:29, 2018.
- [Sha18b] Ohad Shamir. Distribution-specific hardness of learning neural networks. *The Journal of Machine Learning Research*, 19(1):1135–1163, 2018.



- [She08a] Alexander A. Sherstov. Communication lower bounds using dual polynomials. *Bulletin of the EATCS*, 2008.
- [She08b] Alexander A Sherstov. Halfspace matrices. *Computational Complexity*, 17:149–178, 2008.
- [She11] Alexander A Sherstov. The pattern matrix method. *SIAM Journal on Computing*, 40(6):1969–2000, 2011.
- [She12] Alexander A Sherstov. Making polynomials robust to noise. In *Proceedings of the forty-fourth annual ACM symposium on Theory of computing*, pages 747–758, 2012.
- [She13] Alexander A Sherstov. Optimal bounds for sign-representing the intersection of two halfspaces by polynomials. *Combinatorica*, 33(1):73–96, 2013.
- [SSSS10] Shai Shalev-Shwartz, Ohad Shamir, Nathan Srebro, and Karthik Sridharan. Learnability, stability and uniform convergence. *The Journal of Machine Learning Research*, 11:2635–2670, 2010.
- [SSSS17] Shai Shalev-Shwartz, Ohad Shamir, and Shaked Shammah. Failures of gradient-based deep learning. In *Proceedings of the 34th International Conference on Machine Learning-Volume 70*, pages 3067–3075, 2017.
- [ST17] Rocco A Servedio and Li-Yang Tan. What circuit classes can be learned with non-trivial savings? In *8th Innovations in Theoretical Computer Science Conference (ITCS 2017)*. Schloss Dagstuhl-Leibniz-Zentrum fuer Informatik, 2017.
- [STT12] Rocco Servedio, Li-Yang Tan, and Justin Thaler. Attribute-efficient learning and weight-degree tradeoffs for polynomial threshold functions. In *Conference on Learning Theory*, pages 14–1, 2012.

- [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.
- [SW14] Adrien Saumard and Jon A Wellner. Log-concavity and strong log-concavity: a review. *Statistics surveys*, 8:45, 2014.
- [SW19] Alexander A Sherstov and Pei Wu. Near-optimal lower bounds on the threshold degree and sign-rank of  $ac_0$ . In *Proceedings of the 51st Annual ACM SIGACT Symposium on Theory of Computing*, pages 401–412, 2019.
- [SZB21] Min Jae Song, Ilias Zadik, and Joan Bruna. On the cryptographic hardness of learning single periodic neurons. *arXiv preprint arXiv:2106.10744*, 2021.
- [Szö09] Balázs Szörényi. Characterizing statistical query learning: simplified notions and proofs. In *International Conference on Algorithmic Learning Theory*, pages 186–200. Springer, 2009.
- [Tal17] Avishay Tal. Tight bounds on the Fourier spectrum of  $AC^0$ . In *32nd Computational Complexity Conference (CCC 2017)*. Schloss Dagstuhl-Leibniz-Zentrum fuer Informatik, 2017.
- [Tia17] Yuandong Tian. An analytical formula of population gradient for two-layered relu network and its applications in convergence and critical point analysis. In *Proceedings of the 34th International Conference on Machine Learning, ICML 2017*, volume 70, pages 3404–3413. PMLR, 2017.
- [Tie23] Stefan Tiegel. Hardness of agnostically learning halfspaces from worst-case lattice problems. In *The Thirty Sixth Annual Conference on Learning Theory*, pages 3029–3064. PMLR, 2023.

- [TJ<sup>+</sup>16] Florian Tramèr, Fan Zhang 0022, Ari Juels, Michael K. Reiter, and Thomas Ristenpart. Stealing machine learning models via prediction apis. *CoRR*, abs/1609.02943, 2016.
- [Tsy08] A.B. Tsybakov. *Introduction to Nonparametric Estimation*. Springer Series in Statistics. Springer New York, 2008.
- [Vad17] Salil Vadhan. On learning vs. refutation. In *Conference on Learning Theory*, pages 1835–1848. PMLR, 2017.
- [Val84] Leslie G Valiant. A theory of the learnable. *Communications of the ACM*, 27(11):1134–1142, 1984.
- [Vap98] V.N. Vapnik. *Statistical Learning Theory*. A Wiley-Interscience publication. Wiley, 1998.
- [Vap00] Vladimir N. Vapnik. *The nature of statistical learning theory*. Statistics for Engineering and Information Science. Springer-Verlag, New York, second edition, 2000.
- [VdG00] Sara A Van de Geer. *Empirical Processes in M-estimation*, volume 6. Cambridge university press, 2000.
- [Ver18] Roman Vershynin. *High-dimensional probability: An introduction with applications in data science*, volume 47. Cambridge university press, 2018.
- [VRPS21] Gal Vardi, Daniel Reichman, Toniann Pitassi, and Ohad Shamir. Size and depth separation in approximating natural functions with neural networks. *arXiv preprint arXiv:2102.00314*, 2021.
- [VSS<sup>+</sup>22] Kiran Vodrahalli, Rakesh Shivanna, Mahesh Sathiamoorthy, Sagar Jain, and Ed Chi. Algorithms for efficiently learning low-rank neural networks, 2022.

- [Vu98] Van H Vu. On the infeasibility of training neural networks with small mean-squared error. *IEEE Transactions on Information Theory*, 44(7):2892–2900, 1998.
- [Vu06] VH Vu. On the infeasibility of training neural networks with small mean-squared error. *IEEE Transactions on Information Theory*, 44(7):2892–2900, 2006.
- [VW19a] Santosh Vempala and John Wilmes. Gradient descent for one-hidden-layer neural networks: Polynomial convergence and sq lower bounds. In *COLT*, volume 99, 2019.
- [VW19b] Santosh Vempala and John Wilmes. Gradient descent for one-hidden-layer neural networks: Polynomial convergence and sq lower bounds. In *Conference on Learning Theory*, pages 3115–3117, 2019.
- [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.
- [YS20] Gilad Yehudai and Ohad Shamir. Learning a single neuron with gradient methods. *arXiv preprint arXiv:2001.05205*, 2020.
- [ZBH<sup>+</sup>21] Chiyuan Zhang, Samy Bengio, Moritz Hardt, Benjamin Recht, and Oriol Vinyals. Understanding deep learning (still) requires rethinking generalization. *Communications of the ACM*, 64(3):107–115, 2021.
- [Zol84] Vladimir Mikhailovich Zolotarev. Probability metrics. *Theory of Probability & Its Applications*, 28(2):278–302, 1984.
- [ZPS17] Qiuyi Zhang, Rina Panigrahy, and Sushant Sachdeva. Electron-proton dynamics in deep learning. *CoRR*, abs/1702.00458, 2017.

- [ZSJ<sup>+</sup>17a] Kai Zhong, Zhao Song, Prateek Jain, Peter L. Bartlett, and Inderjit S. Dhillon. Recovery guarantees for one-hidden-layer neural networks. In *ICML*, volume 70, pages 4140–4149. JMLR.org, 2017.
- [ZSJ<sup>+</sup>17b] Kai Zhong, Zhao Song, Prateek Jain, Peter L Bartlett, and Inderjit S Dhillon. Recovery guarantees for one-hidden-layer neural networks. In *International conference on machine learning*, pages 4140–4149. PMLR, 2017.
- [ZSWB22] Ilias Zadik, Min Jae Song, Alexander S. Wein, and Joan Bruna. Lattice-based methods surpass sum-of-squares in clustering, 2022.
- [ZYWG19a] Xiao Zhang, Yaodong Yu, Lingxiao Wang, and Quanquan Gu. Learning one-hidden-layer relu networks via gradient descent. In Kamalika Chaudhuri and Masashi Sugiyama, editors, *The 22nd International Conference on Artificial Intelligence and Statistics, AISTATS 2019, 16-18 April 2019, Naha, Okinawa, Japan*, volume 89 of *Proceedings of Machine Learning Research*, pages 1524–1534. PMLR, 2019.
- [ZYWG19b] Xiao Zhang, Yaodong Yu, Lingxiao Wang, and Quanquan Gu. Learning one-hidden-layer relu networks via gradient descent. In *The 22nd International Conference on Artificial Intelligence and Statistics*, pages 1524–1534. PMLR, 2019.