

Weighted Polynomial Approximations: Limits for Learning and Pseudorandomness

Mark Bun* Thomas Steinke†

`{mbun,tsteinke}@seas.harvard.edu`

December 8, 2014

Abstract

Polynomial approximations to boolean functions have led to many positive results in computer science. In particular, polynomial approximations to the sign function underly algorithms for agnostically learning halfspaces, as well as pseudorandom generators for halfspaces. In this work, we investigate the limits of these techniques by proving inapproximability results for the sign function.

Firstly, the “polynomial regression” algorithm of Kalai et al. (SIAM J. Comput. 2008) shows that halfspaces can be learned with respect to log-concave distributions on \mathbb{R}^n in the challenging agnostic learning model. The power of this algorithm relies on the fact that under log-concave distributions, halfspaces can be approximated arbitrarily well by low-degree polynomials. We ask whether this technique can be extended beyond log-concave distributions, and establish a negative result. We show that polynomials of any degree cannot approximate the sign function to within arbitrarily low error for a large class of non-log-concave distributions on the real line, including those with densities proportional to $\exp(-|x|^{0.99})$. This impossibility result extends to multivariate distributions, and thus gives a strong limitation on the power of the polynomial regression algorithm for halfspaces.

Secondly, we investigate the derandomization of Chernoff-type concentration inequalities. Chernoff-type tail bounds on sums of independent random variables have pervasive applications in theoretical computer science. Schmidt et al. (SIAM J. Discrete Math. 1995) showed that these inequalities can be established for sums of random variables with only $O(\log(1/\delta))$ -wise independence, for a tail probability of δ . We show that their results are tight up to constant factors.

These results rely on techniques from weighted approximation theory, which studies how well functions on the real line can be approximated by polynomials under various distributions. We believe that these techniques will have further applications in other areas of theoretical computer science.

*Harvard University, School of Engineering and Applied Sciences. Supported by an NDSEG Fellowship and NSF grant CNS-1237235.

†Harvard University, School of Engineering and Applied Sciences. Supported by NSF grant CCF-1116616 and the Lord Rutherford Memorial Research Fellowship.

1 Introduction

Approximation theory is a classical area of mathematics that studies how well functions can be approximated by simpler ones. It has found many applications in computer science. Most of these applications of approximation theory focus on the approximation of functions by polynomials in the uniform norm (or infinity norm). For instance, *approximate degree*, which captures how well a boolean function can be approximated by low-degree polynomials in the uniform norm, underlies important lower bounds in circuit complexity [Bei93, Bei94, She09], quantum query complexity [BBC⁺01, AS04], and communication complexity [She08]. It also underlies state-of-the-art algorithms in learning theory [KKMS08, KS04], streaming [HNO08], and in spectral methods [SV14].

While it is compelling to study polynomial approximations under the uniform norm, there are scenarios where it is more natural to study *weighted polynomial approximations*, where error is measured in terms of an L_p norm under some distribution. For instance, in agnostic learning, the polynomial regression algorithm of Kalai et al. [KKMS08] has guarantees based on how well functions in a concept class of interest can be approximated by low-degree polynomials in L_1 distance.

In this work, we show how ideas from weighted approximation theory can yield tight lower bounds for several problems in theoretical computer science. As our first application, we establish a strong limitation on the distributions under which halfspaces can be learned using the polynomial regression algorithm of Kalai et al. Second, in the area of derandomization, we give a tight characterization of the amount of k -wise independence necessary to establish Chernoff-like concentration inequalities.

1.1 Agnostically Learning Halfspaces

Halfspaces are a fundamental concept class in machine learning, both in theory and in practice.¹ Their study dates back to the Perceptron algorithm of the 1950s. Halfspaces serve as building blocks in many applications, including boosting and kernel methods.

Halfspaces can be learned in the PAC model [Val84] either by solving a linear program, or via simple iterative update algorithms (e.g. the Perceptron algorithm). However, learning halfspaces with classification noise is a much more difficult problem, and often needs to be dealt with in practice.

In this work, we study a challenging model of *adversarial noise* – the agnostic learning model of Kearns et al. [KSSH94]. In this model, a learner has access to examples drawn from a distribution \mathcal{D} on $X \times \{\pm 1\}$ and must output a hypothesis $h : X \rightarrow \{\pm 1\}$ such that

$$\mathbb{P}_{(x,y) \sim \mathcal{D}} [h(x) \neq y] \leq \text{opt} + \varepsilon,$$

where opt is the error of the best concept in the concept class – that is, $\text{opt} = \min_{f \in \mathcal{C}} \mathbb{P}_{(x,y) \sim \mathcal{D}} [f(x) \neq y]$.

The theory of agnostic learning is not well-understood, even in the case of halfspaces. Positive results for efficient agnostic learning of high-dimensional halfspaces are restricted to limited classes of distributions.² For instance, halfspaces can be learned under the uniform distribution over the

¹A halfspace is a function $f : \mathbb{R}^n \rightarrow \{\pm 1\}$ given by $f(x) = \text{sgn}(w \cdot x - \theta)$ for $w \in \mathbb{R}^n$ and $\theta \in \mathbb{R}$, where $\text{sgn}(x) = 1$ if $x \geq 0$ and $\text{sgn}(x) = -1$ otherwise.

²An efficient algorithm is one which runs in time polynomial in the dimension n for any constant $\varepsilon > 0$ – that is, time $n^{O_\varepsilon(1)}$.

hypercube or the unit sphere, or on any log-concave distribution [KKMS08]. On the negative side, a variety of both computational and information-theoretic hardness results are known. For instance, proper agnostic learning of halfspaces (where the learner is required to output a hypothesis that is itself a halfspace) is known to be NP-hard [FGKP06]. Moreover, agnostically learning halfspaces under arbitrary distributions is as hard as PAC learning DNFs [LBW95], which is a longstanding open problem.

There is essentially only one known technique for agnostically learning high-dimensional halfspaces: the L_1 regression algorithm [KKMS08], which we discuss in more detail in Section 2.2. In its most general form, the algorithm selects a linear space of functions $\mathcal{H} \subset \{h : X \rightarrow \mathbb{R}\}$. After drawing a number of examples (x_i, y_i) from \mathcal{D} , it computes

$$h^* = \operatorname{argmin}_{h \in \mathcal{H}} \sum_i |h(x_i) - y_i|.$$

The output of the algorithm is $\operatorname{sgn}(h^*(x) - t)$ for some t . We need to ensure that the minimisation can be computed efficiently (e.g. by linear programming) and that every concept $f \in \mathcal{C}$ can be approximated by some $h \in \mathcal{H}$ – that is $\mathbb{E}_{x \sim \mathcal{D}} [|h(x) - f(x)|] \leq \varepsilon$. If this is the case, then \mathcal{C} is agnostically learnable in time $\operatorname{poly}(|\mathcal{H}|)$.

Kalai et al. (and most subsequent work on learning using L_1 regression, e.g. [KOS08, GKK08, BOW10, KKM13, FK14]) chose \mathcal{H} to be the class of low-degree polynomials. They showed that under certain classes of distributions, every halfspace can be approximated by a polynomial of degree $O_\varepsilon(1)$, and hence halfspaces are agnostically learnable in time $n^{O_\varepsilon(1)}$.

Distributional assumptions arise because we use an L_1 approximation measure (namely $\mathbb{E}_{x \sim \mathcal{D}} [|h(x) - f(x)|]$), which depends on the distribution. A distribution-independent approximation would require an L_∞ approximation, which is too much to hope for in many circumstances.

1.1.1 Our Results

Can we weaken the distributional assumptions required for learning halfspaces using current techniques? Our result addressing this question (Theorem 2) is a negative one. We show that polynomial approximations to halfspaces do not exist for a large class of distributions, namely:

Definition 1. *An absolutely continuous distribution \mathcal{D} on \mathbb{R} is a log-superlinear (LSL) distribution if there exist $C > 0$ and $\gamma \in (0, 1)$ such that the density w of \mathcal{D} satisfies $w(x) \geq C \exp(-|x|^\gamma)$.³*

Theorem 2. *For any LSL distribution \mathcal{D} , there exists $\varepsilon > 0$ such that no polynomial (of any degree) can approximate the sign function with L_1 error less than ε with respect to \mathcal{D} .*

In particular, this implies that the polynomial regression algorithm is not able to agnostically learn thresholds on the real line to within arbitrarily small error. Note that this result does not rule out the possibility that halfspaces can be agnostically learned by other techniques. Indeed, the classic approach of empirical risk minimization (see [KSSH94] and the references therein) gives an efficient algorithm for learning thresholds (which are halfspaces in one dimension) under arbitrary distributions. Thus the problem of learning real thresholds under LSL distributions is an explicit example for which polynomial regression fails while other techniques can succeed.

³The name log-superlinear comes from the fact that the tails of the probability density function of a LSL distribution are heavier than that of the log-linear Laplace distribution.

If we were to take $\gamma \geq 1$, the probability density function $C(\gamma)e^{-|x|^\gamma}$ (where $C(\gamma)$ is a normalising constant) would give a log-concave distribution, in which case Kalai et al. [KKMS08] show that good polynomial approximations to halfspaces exist. Thus our result gives a threshold between where polynomial approximations to halfspaces exist and where they do not.

Our result for thresholds extends readily to an impossibility result for learning halfspaces over \mathbb{R}^n :

Theorem 3. *For any product distribution \mathcal{D} on \mathbb{R}^n with a LSL marginal distribution on some coordinate, there exists $\varepsilon > 0$ and a halfspace h such that no polynomial can approximate h with L_1 error less than ε with respect to \mathcal{D} .*

Our result echoes prior work establishing the limits of *uniform* polynomial approximations for various concept classes. For instance, the seminal work of Minsky and Papert [MP72] showed that there is an *intersection* of two halfspaces over \mathbb{R}^n which cannot be represented as the sign of any polynomial. Building on work of Nisan and Szegedy [NS94], Paturi [Pat92] gave tight lower bounds for uniform approximations to symmetric boolean functions. This, and subsequent work on lower bounds for approximate degree, immediately imply limitations for distribution-independent agnostic learning via polynomial regression. Klivans and Sherstov [KS10] also showed a strong generalization of Paturi’s result to disjunctions, giving limitations on how well they can be approximated by linear combinations of arbitrary features. By contrast to all of these results, our work shows a strong limitation for certain *distribution-dependent* polynomial approximations.

In the distribution-dependent setting, Feldman and Kothari [FK14] showed that polynomial regression cannot be used to learn disjunctions with respect to symmetric distributions on the hypercube. Recent work of Daniely et al. [DLS14] also uses ideas from approximation theory to show limitations on broad class of regression and kernel-based methods for learning halfspaces, even under a margin assumption. While our results only apply to polynomial regression, they hold for approximations of arbitrarily high complexity (i.e. degree), and for a large class of natural distributions.

The proof of Theorem 2 relies on several Markov-type inequalities for weighted polynomial approximations. These are generalizations of the classical Markov inequality for uniform approximations, which gives a bound on the derivative of a low-degree polynomial that is bounded on the unit interval:

Theorem 4 ([Mar90]). *Let p be a polynomial of degree d with $|p(x)| \leq 1$ on the interval $[-1, 1]$. Then $|p'(x)| \leq d^2$ on $[-1, 1]$.*

Early work on the approximate degree of boolean functions [NS94, Pat92] used Markov’s inequality to get tight lower bounds on the degree of uniform approximations to symmetric functions. For weighted approximations under LSL distributions, we actually get a much stronger statement. It turns out that under LSL distributions, the derivative of a bounded polynomial near the origin is at most a constant *independent of degree*. With this powerful fact in hand, the proof of Theorem 2 is quite simple. Consider the threshold function $f(t) = \text{sgn}(t)$. Since f has a “jump” at zero, any good polynomial approximation to f must be bounded and have a large derivative near zero. The higher quality the approximation, the larger a derivative we need. But since the derivative of any polynomial is bounded by a constant, we cannot get arbitrarily good approximations to f using polynomials.

We give the full proof in Section 2.4, and discuss the multivariate generalization in Section 2.5.

1.1.2 Related Work

There is a rich literature on lower bounds for agnostic learning. In the case of *proper* agnostic learning Feldman et. al [FGKP06] gave an optimal NP-hardness result for even weakly agnostically learn halfspaces over \mathbb{Q}^n . Guruswami and Raghavendra [GR06] showed that the same is true even for halfspaces on the boolean hypercube.

There has also been a line of work giving representation-independent hardness of learning halfspaces based on cryptographic assumptions. Feldman et. al [FGKP06] and Klivans and Sherstov [KS09] showed that, assuming the security of certain public key encryption schemes, it is hard to even PAC learn thresholds and intersections of halfspaces, respectively. These results imply that it is hard to agnostically learn a single halfspace in the harsh noise regime, i.e. when opt is very close to $\frac{1}{2}$. Shalev-Schwartz et al. [SSSS11] further showed that halfspaces cannot be efficiently learned even under a large margin assumptions.

There has also been extensive work proving unconditional lower bounds for restricted learning algorithms. One well-studied restriction on a learner is that it operates in Kearns' statistical query (SQ) model [Kea98, BFJ⁺94, KS07, FLS11]. This model captures L_1 regression, as well as essentially every technique known for learning (besides Gaussian elimination). Very recently, Dachman-Soled et al. [DFT⁺14] showed that polynomial regression is in fact essentially the optimal SQ algorithm for agnostic learning with respect to product distributions on the hypercube.

The limitations we prove for polynomial regression do not rule out the existence of other agnostic learning algorithms, including those using L_1 regression with different feature spaces. Wimmer [Wim10] showed how to use a different family of basis functions to learn halfspaces over symmetric distributions on the hypercube. Subsequent work of Feldman and Kothari [FK14] improved the running time in the special case of disjunctions. We leave it as an intriguing open question to determine whether other basis functions can be used to learn halfspaces under LSL distributions.

1.2 Tail Bounds for Limited Independence

The famous Hoeffding bound [Hoe63] implies that if $X \in \{\pm 1\}^n$ is a uniform random variable and $r \in \mathbb{R}^n$ is fixed, then, for all $T \geq 0$,

$$\mathbb{P}_X[|X \cdot r| \geq T] \leq 2e^{-\frac{T^2}{2\|r\|_2^2}}.$$

We ask the following question:

For what *pseudorandom* X is the Hoeffding bound true?

More precisely, given T and δ , can we construct a pseudorandom X such that $\mathbb{P}_X[|X \cdot r| \geq T] \leq \delta$ for all $r \in \{\pm 1\}^n$?⁴ Of particular interest is the parameter regime $\delta = 1/\text{poly}(n)$ and $T = \Theta(\|r\|_2 \sqrt{\log(1/\delta)})$. The probabilistic method gives a non-constructive proof that there exists such an X which can be sampled with seed length $O(\log(n/\delta))$. The challenge is to give an explicit construction of such an X which can be *efficiently* sampled with a short seed.

This is a very natural pseudorandomness question: Concentration of measure is a fundamental property of independent random variables and one of the key objectives of pseudorandomness research is to replicate such properties for variables with low entropy. Finding a pseudorandom X

⁴For simplicity we restrict our attention to $r \in \{\pm 1\}^n$.

exhibiting good concentration is also a relaxation of a more general and well-studied pseudorandomness question, namely constructing pseudorandom generators that fool linear threshold functions [DGJ⁺09, MZ10, GOWZ10, DSTW10]. This can also be viewed as a special case of constructing pseudorandom generators for space-bounded computation [Nis92, INW94, Rei08, BRRY10, BV10, KNP11, RSV13].

For $\delta = 1/\text{poly}(n)$ and $T = \Theta(\sqrt{n \log(1/\delta)})$, we can construct generators X with seed length $O(\log^2 n)$ using a variety of methods (including [Nis92, MZ10]). In particular, it suffices for X to be $O(\log(1/\delta))$ -wise independent:

Theorem 5 (Tail Bound for Limited Independence). *Let $n \geq 1$, $\eta > 0$, and $\delta \in (0, 1)$ be given. Let $X \in \{\pm 1\}^n$ be k -wise independent for $k = 2\lceil \eta \log_e(1/\delta) \rceil$. Set $T = e^{(\eta+1)/2\eta} \sqrt{k} \|r\|_2$. Then, for all $r \in \mathbb{R}^n$,*

$$\mathbb{P}[|X \cdot r| \geq T] \leq \delta.$$

A k -wise independent $X \in \{\pm 1\}^n$ can be sampled with seed length $O(k \cdot \log n)$ [ABI85]. Another construction which achieves seed length $O(\log n \cdot \log(1/\delta))$ is to sample X from a small-bias space [NN93]. Very recently, Gopalan et al. [GKM14] constructed a new generator with seed length $\tilde{O}(\log(n/\delta))$, which is nearly optimal.

In this work, we ask whether the tail bound of Theorem 5 for k -wise independence is tight. That is, can we prove stronger tail bounds for k -wise independent X ?

Question 6. *How much independence is needed for X to satisfy a Hoeffding-like tail bound? That is, what is the minimum $k = k(n, \delta, T)$ for which any k -wise independent $X \in \{\pm 1\}^n$ satisfies*

$$\mathbb{P}_X[|X \cdot r| \geq T] \leq \delta$$

for all $r \in \{-1, 1\}^n$, where \cdot denotes the inner product.

1.2.1 Our Results

Theorem 5 shows that $k(n, \delta, T) \leq O(\log(1/\delta))$ for $T = O(\sqrt{n \log(1/\delta)})$. In this work, we show that this is essentially tight:

Theorem 7. *For $T = c\sqrt{n \log(1/\delta)}$ ($c > 5$), we have $k(n, \delta, T) = \Omega(\log_c(1/\delta))$ for sufficiently large n .*

The only previous lower bound was

$$k(n, \delta, T) \geq \Omega\left(\frac{\log(1/\delta)}{\log n}\right),$$

which holds for any $T \leq n$ and is due to [SSS95]. This is useful if $\delta < n^{-\omega(1)}$, but the lower bound is constant in our parameter regime. This lower bound follows from the fact that a random variable X with support size s cannot give a tail bound with $\delta < 1/s$, and that there exist k -wise independent distributions with support size $s \leq O(n^k)$.

The most natural way to prove Theorem 7 would be to construct a family of k -wise independent distributions that do not satisfy the required tail bound. However, we instead study the *dual* formulation of the problem (following [Baz09, DETT10]) and then use lower bound techniques

from approximation theory. To the best of our knowledge, this indirect approach is novel. Our results imply the existence of k -wise independent distributions with poor tail bounds, but give no immediate indication as to how to construct them!

We now describe the proof idea in slightly more detail. The answer to Question 6 can be posed in terms of the value of a certain linear program. The variables represent the probability distribution of the random variable X and the constraints force X to be k -wise independent. The objective of the linear program is maximize $\mathbb{P}[|X \cdot r| \geq T]$. Thus, the value of the program is at most δ if and only if $k \geq k(n, \delta, T)$. Taking the dual of this linear program and appealing to strong duality yields an alternative characterization of $k(n, \delta, T)$. Namely, $k(n, \delta, T)$ is the smallest k for which the threshold function $F_T(x) = \mathbb{1}(|x| \geq T)$ admits an *upper sandwiching polynomial* of degree k and expectation at most δ . Here, an upper sandwiching polynomial is simply a polynomial p for which $p(x) \geq F_T(x)$ pointwise.

We then use ideas from weighted approximation theory to give a lower bound on k for which such sandwiching polynomials exist. In order to apply these ideas, we make a few symmetrization and approximation arguments to reduce the problem to a continuous one-dimensional problem: Find a degree lower bound for a univariate polynomial that is a good upper sandwich for the function $f_T(x) = \text{sgn}(|x| - T)$, with respect to a Gaussian distribution. As in our proof of Theorem 2, the solution of this problem appeals to a weighted Markov-type inequality. Again, the idea is that an upper sandwich for f_T must have a large jump at the threshold T , which is impossible for low-degree polynomials. The formal proof of this claim is based on a variant of an “infinite-finite range” inequality, which asserts that the weighted norm of a polynomial on the real line is bounded by its norm on a finite interval.

2 Agnostically Learning Halfspaces

The class of log-concave distributions over \mathbb{R}^n (defined below) is essentially the broadest under which we know how to agnostically learn halfspaces. While many distributions used in machine learning are log-concave, such as the normal, Laplace, beta, and Dirichlet distributions, log-concave distributions do not capture everything. For instance, the log-normal distribution and heavier-tailed exponential power law distributions are not log-concave. The main motivating question for this section is whether we can relax the assumption of log-concavity for agnostically learning halfspaces. To this end, we show a negative result: for LSL distributions, agnostic learning of halfspaces will require new techniques.

2.1 Background

Our starting point is the work of Kalai et al. [KKMS08]. Among their results is the following.

Theorem 8 ([KKMS08]). *The concept class of halfspaces over \mathbb{R}^n is agnostically learnable in time $\text{poly}(n^{O_\varepsilon(1)})$ under log-concave distributions.*

A log-concave distribution is an absolutely continuous probability distribution such that the logarithm of the probability density function is concave. For example, the standard multivariate Gaussian distribution on \mathbb{R}^n has the probability density function $x \mapsto e^{-\|x\|_2^2/2}/(2\pi)^{n/2}$. The natural logarithm of this is $-\|x\|_2^2/2 - n/2 \cdot \log(2\pi)$, which is concave. The class of log-concave distributions also includes the Laplace distribution and other natural distributions. However, it

does not contain heavy-tailed distributions (such as power laws) nor non-smooth distributions (such as discrete probability distributions).

Kalai et al. also show that we can agnostically learn halfspaces under the uniform distribution over the hypercube $\{\pm 1\}^n$ or over the unit sphere $\{x \in \mathbb{R}^n : \|x\|_2 = 1\}$.

2.2 The L_1 Regression Algorithm

The results of Kalai et al. are based on the so-called L_1 regression algorithm, which relies on being able to approximate the concept class in question by a low-degree polynomial:

Theorem 9 ([KKMS08]). *Fix a distribution \mathcal{D} on $X \times \{\pm 1\}$ and a concept class $\mathcal{C} \subset \{f : X \rightarrow \{\pm 1\}\}$.⁵ Suppose that, for all $f \in \mathcal{C}$, there exists a polynomial $p : X \rightarrow \mathbb{R}$ of degree at most d such that $\mathbb{E}_{x \sim \mathcal{D}_X} [|p(x) - f(x)|] \leq \varepsilon$, where \mathcal{D}_X is the marginal distribution of \mathcal{D} on X . Then, with probability $1 - \delta$ the L_1 regression algorithm outputs a hypothesis h such that*

$$\mathbb{P}_{(x,y) \sim \mathcal{D}} [h(x) \neq y] \leq \min_{f \in \mathcal{C}} \mathbb{P}_{(x,y) \sim \mathcal{D}} [f(x) \neq y] + \varepsilon$$

in time $\text{poly}(n^d, 1/\varepsilon, \log(1/\delta))$ with access only to examples drawn from \mathcal{D} .

The L_1 regression algorithm solves a linear program to find a polynomial p of degree at most d that minimises $\sum_i |p(x_i) - y_i|$, where (x_i, y_i) are the examples sampled from \mathcal{D} . The hypothesis is then $h(x) = \text{sgn}(p(x) - t)$, where $t \in [-1, 1]$ is chosen to minimise the error of h on the examples.

Given Theorem 9, proving Theorem 8 reduces to showing that halfspaces can be approximated by low-degree polynomials under the distributions we are interested in. It is important to note that making assumptions on the distribution is necessary (barring a major breakthrough): Agnostically learning halfspaces under arbitrary distributions is at least as hard as PAC learning DNF formulas [LBW95]. Moreover, proper learning of halfspaces under arbitrary distributions is known to be NP-hard [FGKP06].

In fact, we can reduce the task of approximating a halfspace to a one-dimensional problem. A halfspace is given by $f(x) = \text{sgn}(w \cdot x - \theta)$ for some $w \in \mathbb{R}^n$ and $\theta \in \mathbb{R}$. It suffices to find a univariate polynomial p of degree at most d such that $\mathbb{E}_{x \sim \mathcal{D}_{w,\theta}} [|p(x) - \text{sgn}(x)|] \leq \varepsilon$, where $\mathcal{D}_{w,\theta}$ is the distribution of $w \cdot x - \theta$ when x is drawn from \mathcal{D}_X . If \mathcal{D}_X is log-concave, then so is $\mathcal{D}_{w,\theta}$.

2.3 On the Density of Polynomials

In this section, we give some intuition for why one might expect that polynomial approximations do not suffice for learning under LSL distributions. It turns out that under a LSL distribution w , polynomials actually fail to be dense in the space $C_0[w]$ of continuous functions vanishing at infinity when weighted by w . This is in stark contrast to the classical Weierstrass approximation theorem, which asserts that the polynomials are dense in C_0 under the uniform weight. These kinds of results address *Bernstein's approximation problem* [Ber24], a precise statement of which is as follows.

⁵Here $X = \mathbb{R}^n$.

Question 10. Let $w : \mathbb{R} \rightarrow [0, 1]$ be a measurable function. Let $C_0[w]$ denote the space of continuous functions f for which $\lim_{|x| \rightarrow \infty} f(x)w(x) = 0$. Under what conditions on w is it true that for every $f \in C_0[w]$, there is a sequence of polynomials $\{p_n\}_{n=1}^\infty$ for which

$$\lim_{n \rightarrow \infty} \|(p_n - f)w\|_\infty = 0?$$

(The choice of the L_∞ norm here appears to make very little difference). If Bernstein's problem admits a positive resolution, we say that the polynomials are *dense* in $C_0[w]$. The excellent survey of Lubinsky [Lub07] presents a number of criteria for when polynomials are dense. The one that is most readily applied was proved by Carleson [Car51] (but appears to be implicit in [IK37]):

Theorem 11. Let w be even and positive with $\log(w(e^x))$ concave. Then the polynomials are dense in $C_0[w]$ iff

$$\int_0^\infty \frac{\log w(x)}{1+x^2} dx = -\infty.$$

This immediately yields the following dichotomy result for exponential power distributions:

Corollary 12. For $\gamma > 0$ and $w_\gamma(x) = \exp(-|x|^\gamma)$, the polynomials are dense in $C_0[w_\gamma]$ iff $\gamma \geq 1$.

In particular, this justifies our assertion that the polynomials fail to be dense in the continuous functions under LSL distributions.

So what does this have to do with agnostically learning halfspaces? Recall that the analysis of the L_1 -regression algorithm of Kalai et al. [KKMS08] reduces approximating a halfspace under a distribution \mathcal{D} to the problem of approximating each threshold function $\text{sgn}(x - \theta)$ under each marginal distribution of \mathcal{D} . So for the algorithm to work, we require \mathcal{D} to have marginals w under which $\text{sgn}(x - \theta)$ can be approximated arbitrarily well by polynomials. Now if the polynomials are dense in $C_0[w]$, then threshold functions can also be approximated arbitrarily well (since $C_0[w]$ is in turn dense in $L_1[w]$). Such an appeal to density actually underlies Kalai et al.'s proof of approximability under log-concave distributions. On the other hand, if the polynomials fail to be dense, then one might conjecture that thresholds cannot be arbitrarily well approximated.

Our result, presented in the next section, confirms the conjecture that even the *sign* function cannot be approximated arbitrarily well by polynomials under LSL distributions.

2.4 Lower Bound for One Variable

Consider the LSL density function

$$w_\gamma(x) := C(\gamma) \exp(-|x|^\gamma)$$

on the reals for $\gamma \in (0, 1)$, where $C(\gamma)$ is a normalizing constant. Define the sign function $\text{sgn}(x) = 1$ if $x \geq 0$ and $\text{sgn}(x) = -1$ otherwise. In this section, we show that for sufficiently small ε , the sign function does not have an L_1 approximation under the distribution w_γ . More formally,

Proposition 13. For any $\gamma \in (0, 1)$, there exists an $\varepsilon = \varepsilon(\gamma)$ such that for any polynomial p ,

$$\int_{\mathbb{R}} |p(x) - \text{sgn}(x)| w_\gamma(x) dx > \varepsilon.$$

The proof is based on the following Markov-type inequality, which roughly says that a bounded polynomial cannot have a large derivative (under the weight w_γ). This implies the claim, since the sign function we are trying to approximate has a large “jump” at the origin.

Lemma 14. *For $\gamma \in (0, 1)$ there is a constant $M(\gamma)$ such that*

$$\sup_{x \in \mathbb{R}} (|p'(x)|w_\gamma(x)) \leq M(\gamma) \int_{\mathbb{R}} |p(x)|w_\gamma(x) dx.$$

Proof. The lemma is a combination of a Markov-type inequality and a Nikolskii-type, available in a survey of Nevai [Nev86]:

Theorem 15 ([NT86], [Nev86, Theorem 4.17.4]). *There exists a constant $C_1(\gamma)$ such that for any polynomial p ,*

$$\int_{\mathbb{R}} |p'(x)|w_\gamma(x) dx \leq C_1(\gamma) \int_{\mathbb{R}} |p(x)|w_\gamma(x) dx.$$

Theorem 16 ([NT87], [Nev86, Theorem 4.17.5]). *There exists a constant $C_2(\gamma)$ such that for any polynomial p ,*

$$\sup_x (|p(x)|w_\gamma(x)) \leq C_2(\gamma) \int_{\mathbb{R}} |p(x)|w_\gamma(x) dx.$$

□

Proof of Proposition 13. Fix $\varepsilon \in (0, 1)$ and suppose p is a polynomial satisfying

$$\int_{\mathbb{R}} |p(x) - \text{sgn}(x)|w_\gamma(x) dx \leq \varepsilon.$$

Since the absolute value of the sign function integrates to 1, this forces

$$\int_{\mathbb{R}} |p(x)|w_\gamma(x) dx \leq 1 + \varepsilon \leq 2.$$

Therefore, we have by Lemma 14 that $|p'(x)|w_\gamma(x) \leq 2M(\gamma)$ for every x .

The idea is now to show that there is some x_0 for which $|p'(x_0)|w_\gamma(x_0) \geq \Omega(1/\varepsilon)$. To see this, let $\delta = 4\varepsilon/C(\gamma)$ and observe that there must exist some $x_+ \in [0, \delta]$ such that $p(x_+) \geq 1/2$. If this were not the case, then we would have

$$\int_{\mathbb{R}} |p(x)|w_\gamma(x) dx \geq \frac{1}{2} \int_0^\delta C(\gamma) \exp(-\delta^\gamma) \geq \varepsilon$$

for δ small enough to make $\exp(-\delta^\gamma) \geq 1/2$, yielding a contradiction. A similar argument shows that there is some $x_- \in [-\delta, 0]$ with $p(x_-) \leq -1/2$. Therefore, by the mean value theorem, there is some $x_0 \in [x_-, x_+]$ with $p'(x_0) \geq 1/2\delta = C(\gamma)/8\varepsilon$. Moreover, because we took δ small enough, we also have $p'(x_0)w_\gamma(x_0) \geq C(\gamma)/16\varepsilon$. This shows that no polynomial ε -approximates sgn as long as $\varepsilon < C/32M$. □

Moreover, the proposition shows that it is impossible to get arbitrarily close polynomial approximations to halfspaces under densities w for which there are constants C and $\gamma \in (0, 1)$ with $w(x) \geq C \exp(-|x|^\gamma)$ for all $x \in \mathbb{R}$. This shows that LSL distributions on \mathbb{R} do not support polynomial approximations to halfspaces.

2.5 Extending the Lower Bound to Multivariate Distributions

It is straightforward to extend the lower bound from the previous section to product distributions with LSL marginals.

Theorem 17. *Let $X = (X_1, \dots, X_n)$ be a random variable over \mathbb{R}^n with density $f_X(x) = w(x_1)f(x_2, \dots, x_n)$. Suppose the density w specifies a univariate γ -LSL distribution. Then there exists an $\varepsilon = \varepsilon(\gamma)$ such that for any polynomial p ,*

$$\int_{\mathbb{R}^n} |p(x_1, \dots, x_n) - \text{sgn}(x_1)| f_X(x_1, \dots, x_n) dx_1 dx_2 \dots dx_n > \varepsilon.$$

That is, the linear threshold function $\text{sgn}(x_1)$ cannot be approximated arbitrarily well by polynomials.

Proof. Let $p(x_1, \dots, x_n)$ be a polynomial, and define a univariate polynomial q by “averaging out” the variables x_2, \dots, x_n :

$$q(x_1) := \int_{\mathbb{R}^{n-1}} p(x_1, \dots, x_n) f(x_2, \dots, x_n) dx_2 \dots dx_n.$$

Then we have

$$\begin{aligned} \int_{\mathbb{R}} |q(x_1) - \text{sgn}(x_1)| w(x_1) dx_1 &= \int_{\mathbb{R}} \left| \int_{\mathbb{R}^{n-1}} (p(x_1, \dots, x_n) - \text{sgn}(x_1)) f(x_2, \dots, x_n) dx_2 \dots dx_n \right| w(x_1) dx_1 \\ &\leq \int_{\mathbb{R}} \left(\int_{\mathbb{R}^{n-1}} |p(x_1, \dots, x_n) - \text{sgn}(x_1)| f(x_2, \dots, x_n) dx_2 \dots dx_n \right) w(x_1) dx_1 \\ &= \int_{\mathbb{R}^n} |p(x_1, \dots, x_n) - \text{sgn}(x_1)| f_X(x_1, \dots, x_n) dx_1 dx_2 \dots dx_n. \end{aligned}$$

By Proposition 13, the latter quantity must be at least $\varepsilon(\gamma)$. □

Let $w_\gamma^n(x) \propto \exp(-(|x_1|^\gamma + \dots + |x_n|^\gamma))$ denote the density of the prototypical multivariate LSL distribution, with each marginal having the same exponential power law distribution. Our impossibility result holds uniformly for every distribution in the sequence $\{w_\gamma^n\}$. That is, for every $\gamma \in (0, 1)$, there exists $\varepsilon = \varepsilon(\gamma)$ for which halfspaces cannot be learned by polynomials under any of the distributions specified by $\{w_\gamma^n\}$.

As a consequence, we get inapproximability results for several natural classes of distributions that dominate $\{w_\gamma^n\}$ by constant factors (i.e. not growing with n).

1. Any power-law distribution, i.e. a distribution with density $\propto \|x\|^{-M}$ for some constant M , since such a distribution dominates every w_γ^n .
2. Multivariate generalizations of the log-normal distribution, i.e. any distribution with density $\propto \exp(-\text{polylog}(\|x\|))$.
3. Multivariate exponential power distributions, which have densities $\propto \exp(-\|x\|^\gamma)$ for $\gamma \in (0, 1)$. These distributions dominate the prototypical w_γ^n by the inequality of ℓ_p -norms:

$$\|x\|^\gamma \leq |x_1|^\gamma + \dots + |x_n|^\gamma$$

for every $0 \leq \gamma \leq 2$.

3 Tail Bounds for Limited Independence

Our proof consists of three steps:

- §3.1 First we reformulate the question of tail bounds for k -wise independent distributions using linear programming duality and symmetrisation. This reduces the problem to proving a degree lower bound on univariate polynomials. Namely we need to give a lower bound on the degree of a polynomial $p : \{0, 1, \dots, n\} \rightarrow \mathbb{R}$ such that $p(i) \geq 0$ for all i , $p(i) \geq 1$ if $|i - n/2| \geq T$, and $\mathbb{E}[p(i)] \leq \delta$, where i is drawn from the binomial distribution.
- §3.2 We then transform the problem from one about polynomials with a discrete domain to one about polynomials with a continuous domain. This amounts to showing that, since $\mathbb{E}[p(i)] \leq \delta$ with respect to the binomial distribution, we can bound $\mathbb{E}[p(x + n/2)]$ with respect to a truncated Gaussian distribution on x .
- §3.3 Finally we can apply the tools of weighted approximation theory. We know that $p(x + n/2)$ is small for x near the origin, but $p(T + n/2) \geq 1$. We show that any low-degree polynomial that is bounded near the origin cannot grow too quickly. This implies that p must have high degree.

3.1 Dual Formulation

Question 6 from the introduction is equivalent to finding the smallest k for which the value of the following linear program is at most δ .

Linear Program Formulation of Question 6

$$\begin{aligned}
 \max_{\psi} \quad & \sum_{x \in \{-1, 1\}^n} \psi(x) F_T(x) \\
 \text{s.t.} \quad & \sum_{x \in \{-1, 1\}^n} \psi(x) \chi_S(x) = 0 && \text{for all } |S| \leq k \\
 & \sum_{x \in \{-1, 1\}^n} \psi(x) = 1 \\
 & 0 \leq \psi(x) \leq 1 && \text{for all } x \in \{-1, 1\}^n.
 \end{aligned}$$

Here, $F_T(x) = 1$ if $|x| \geq T$ and is 0 otherwise, and $\chi_S(x)$ is the Fourier character corresponding to $S \subseteq [n]$.

If we set $\mathbb{P}_X[X = x] = \phi(x)$, then the constraints impose that X is a k -wise independent distribution, while the objective function is $\mathbb{P}_X\left[\left|\sum_{i \in [n]} X_i\right| \geq T\right]$. Thus the above linear program finds the k -wise independent distribution with the worst tail bound. If the value of the program is at most δ , then all k -wise independent distributions satisfy the tail bound, as required.

Taking the dual of the above linear program yields the following.

Dual Formulation of Question 6

$$\begin{aligned}
& \min_p 2^{-n} \sum_{x \in \{-1, 1\}^n} p(x) \\
& \text{s.t. } \deg(p) \leq k \\
& \quad p(x) \geq F_T(x) \quad \text{for all } x \in \{-1, 1\}^n.
\end{aligned}$$

By strong duality, the value of the dual linear program is the same as that of the primal.

The multilinear polynomial p as an “upper sandwich” of F_T – that is, $p \geq F_T$ and $\mathbb{E}_{X \in \{\pm 1\}^n} [p(X)]$ is minimal. Therefore, $k(n, \delta, T)$ is the smallest k for which F_T admits an upper sandwiching polynomial of degree k with expectation δ .

Consider the shifted univariate symmetrization of F_T

$$F'_T(x) = \begin{cases} 1 & \text{if } |x - n/2| \geq T \\ 0 & \text{otherwise.} \end{cases}$$

By applying the well-known Minsky-Papert symmetrization [MP72] to the dual formulation above, we get the following characterization.

Theorem 18. *The quantity $k(n, \delta, T)$ from Question 6 is the smallest k for which there exists a degree- k univariate polynomial $p : \{0, \dots, n\} \rightarrow \mathbb{R}$ such that*

1. $p(i) \geq F'_T(i)$ for all $0 \leq i \leq n$ and
2. $2^{-n} \sum_{i=0}^n \binom{n}{i} p(i) \leq \delta$.

The upper bound on $k(n, \delta, T)$ (Theorem 5) is proved (in the appendix) by showing that

$$p(i) = \left(\frac{i - n/2}{T} \right)^k$$

satisfies the requirements of Theorem 18 for an appropriate even k .⁶ So this characterisation does in fact capture how upper bounds are proved. The fact that it is a tight characterisation allows us to prove that a barrier to the technique is in fact an impossibility result.

With this characterisation of our problem, we may move on to proving inapproximability results.

3.2 A Continuous Version

To apply techniques from the theory of weighted polynomial approximations, we move to polynomials on a continuous domain. We replace the binomial distribution upon which Theorem 18 evaluates p with a Gaussian distribution.

Define the probability density function

$$w(x) = \frac{1}{\sqrt{\pi}} e^{-x^2}.$$

We define the L_∞ norm with respect to the weight w :

$$\|g\|_{L_\infty(S)} = \sup_{x \in S} |g(x)| w(x).$$

Now we can give the continuous version of the problem:

⁶While our results show that this polynomial is *asymptotically* optimal, numerical experiments have shown that it is not exactly optimal.

Theorem 19. Let $T = c\sqrt{n \log(1/\delta)}$ for $c \geq 5$, and $d = k(n, \delta, T)$. Assume $n \geq (12c)^2(3 \log(1/\delta))^3$. Then for $T' = 4cT/\sqrt{n}$, there is a degree d polynomial q such that

1. $q(T') = q(-T') \geq 1$ and
2. $\|q\|_{L_\infty[-\sqrt{d}, \sqrt{d}]} \leq \delta^{0.9}(n+1)$.

The following lemma is key to moving from the discrete to the continuous setting. It shows that if a polynomial is bounded at evenly spaced points, then it must also be bounded between those points, assuming the number of points is sufficiently large relative to the degree.

Lemma 20. [EZ64, RC66, NS94] Let q be a polynomial of degree d such that $|q(i)| \leq 1$ for $i = 0, 1, \dots, m$, where $3d^2 \leq m$. Then $|q(x)| \leq \frac{3}{2}$ for all $x \in [0, m]$.

Proof. Let $a = \max_{x \in [0, m]} |q'(x)|$. Then by the mean value theorem, $|q(x)| \leq 1 + a/2$ for $x \in [0, m]$. By Markov's inequality ([Mar90], see also [Che82]),

$$a \leq \frac{2d^2(1 + a/2)}{m}.$$

Rearranging gives

$$\frac{a}{2 + a} \leq \frac{d^2}{m} \leq \frac{1}{3}.$$

Therefore, $a \leq 1$, and hence $|q(x)| \leq \frac{3}{2}$ for $x \in [0, m]$. \square

We also require the following anti-concentration lemma.

Lemma 21.

$$\binom{n}{n/2 + \alpha\sqrt{n}} \geq \frac{2^{n-6\alpha^2}}{n+1}.$$

Proof. It is well known via Stirling's approximation that $\binom{n}{k} \geq 2^{nH(k/n)}/(n+1)$, where $H(\cdot)$ denotes the binary entropy function. We estimate

$$\begin{aligned} H\left(\frac{1}{2} + \frac{\alpha}{\sqrt{n}}\right) &\geq \left(\frac{1}{2} + \frac{\alpha}{\sqrt{n}}\right) \left(1 - \frac{2\alpha}{(\log 2)\sqrt{n}}\right) + \left(\frac{1}{2} - \frac{\alpha}{\sqrt{n}}\right) \left(1 + \frac{2\alpha}{(\log 2)\sqrt{n}}\right) \\ &\geq 1 - \frac{4\alpha^2}{(\log 2)n}, \end{aligned}$$

which concludes the proof. \square

Proof of Theorem 19. Let p be the polynomial promised by Theorem 18. By Theorem 5, we know that $d \leq 3 \log(1/\delta)$. Define

$$q(x) = p(x\sqrt{n}/4c + n/2).$$

Then $q(\pm T') = p(\pm T + n/2) \geq F_T(\pm T + n/2) = 1$, dispensing with the first claim.

Now for all integers i in the interval $n/2 \pm \sqrt{nd}/4c$, we have

$$2^{-n} \binom{n}{i} |p(i)| \leq \delta$$

and hence, by Lemma 21,

$$|p(i)| \leq \frac{2^n \delta}{\binom{n}{n/2 + \sqrt{nd}/4c}} \leq (n+1)\delta 2^{6d/16c^2} \leq (n+1)\delta^{1-18/16c^2} \leq (n+1)\delta^{0.9}.$$

By Lemma 20, $|p(x)| \leq \frac{3}{2}(n+1)\delta^{0.9}$ on the whole interval $n/2 \pm \sqrt{nd}/4c$. Thus $|q(x)| \leq \frac{3}{2}(n+1)\delta^{0.9}$ on $[-\sqrt{d}, \sqrt{d}]$, completing the proof. \square

3.3 The Lower Bound

Now we state the result we need from approximation theory. The following ‘‘infinite-finite range inequality’’ shows that the norm of weighted polynomial on the real line is determined by its norm on a finite interval around the origin. Thus, an upper bound on the magnitude of a polynomial near the origin yields a bound on its growth away from the origin.. We will apply this to the polynomial given to us in Theorem 19.

Theorem 22. *For any polynomial p of degree d and $B > 1$,*

$$\|p\|_{L_\infty(\mathbb{R} \setminus [-B\sqrt{d}, B\sqrt{d}])} \leq (2eB)^d \exp(-B^2 d) \|p\|_{L_\infty[-\sqrt{d}, \sqrt{d}]}.$$

The proof follows [Lub07, Theorem 6.1] and [Nev86, Theorem 4.16.12].

Proof. Let \tilde{p} be a polynomial of degree d . Let $T_d(x)$ denote the d th Chebyshev polynomial of the first kind [Che82]. By the extremal properties of T_d , we have

$$|\tilde{p}(x)| \leq |T_d(x)| \left(\max_{t \in [-1, 1]} |\tilde{p}(t)| \right) \leq (2|x|)^d \left(\max_{t \in [-1, 1]} |\tilde{p}(t)| \right)$$

for $|x| \geq 1$. Rescaling $p(x) = \tilde{p}(x/\sqrt{d})$ yields

$$|p(x)| \leq \left(\frac{2|x|}{\sqrt{d}} \right)^d \left(\max_{t \in [-\sqrt{d}, \sqrt{d}]} |p(t)| \right) \leq \sqrt{\pi} e^d \left(\frac{2|x|}{\sqrt{d}} \right)^d \|p\|_{L_\infty[-\sqrt{d}, \sqrt{d}]}$$

for $|x| \geq \sqrt{d}$. Now let $|x| = B\sqrt{d}$ for some $B > 1$. Then

$$|p(x)| w(x) \leq e^d (2B)^d \exp(-B^2 d) \|p\|_{L_\infty[-\sqrt{d}, \sqrt{d}]}.$$

Since the coefficient $(2eB)^d \exp(-B^2 d)$ is decreasing in B , this proves the claim. \square

The above approximation theory result, combined with our continuous formulation Theorem 19, enables us to complete the proof.

Theorem 23. *Let $T = c\sqrt{n \log(1/\delta)}$ for $c \geq 5$. Assume $n \geq (12c)^2 (3 \log(1/\delta))^3$. Then $k(n, \delta, T) > \log(1/\delta)/9 \log c$.*

Proof. Let q be the polynomial given by Theorem 19. Let $T' = 4cT/\sqrt{n}$, $d = \log(1/\delta)/9 \log c$, and $B = T'/\sqrt{d} = 12c^2 \sqrt{\log c}$. For the sake of contradiction, we suppose that q satisfies the conditions of Theorem 19, but $\deg(q) \leq d$. Then

$$\|q\|_{L_\infty(\mathbb{R} \setminus [-B\sqrt{d}, B\sqrt{d}])} = \|q\|_{L_\infty(\mathbb{R} \setminus [-T', T'])} \geq \frac{\exp(-T'^2)}{\sqrt{\pi}}.$$

On the other hand, applying Theorem 22, gives

$$\|q\|_{L_\infty(\mathbb{R} \setminus [-B\sqrt{d}, B\sqrt{d}])} \leq (2eB)^d \exp(-T'^2) \delta^{0.9} (n+1).$$

Combining the two inequalities gives

$$\frac{1}{\sqrt{\pi}} \leq (2eB)^d \delta^{0.9} (n+1) \leq \left(24ec^2 \sqrt{\log(c)}\right)^{\log(1/\delta)/9 \log(c)} \delta^{0.9} (n+1) \leq \delta^{1/3} (n+1),$$

which is a contradiction. □

Theorem 23 yields Theorem 7.

4 Further Work

Our negative results naturally suggest a number of directions for future work.

Are halfspaces agnostically learnable under LSL distributions? Our negative result does not even necessarily rule out the use of L_1 regression for this task: The polynomial regression algorithm of Kalai et al. [KKMS08] is in fact quite flexible. Nothing is really special about the basis of low-degree monomials, and the algorithm works equally well over any small, efficiently evaluable “feature space”. That is, if we can show that halfspaces are well-approximated by linear combinations of features from a feature space \mathcal{F} under a distribution \mathcal{D} , then we can agnostically learn halfspaces with respect to \mathcal{D} in time proportional to $|\mathcal{F}|$. Could one hope for such approximations? Wimmer [Wim10] and Feldman and Kothari [FK14] have shown how to use non-polynomial basis functions to obtain faster learning algorithms on the boolean hypercube. On the other hand, recent work of Dachman-Soled et al. [DFT⁺14] shows that, at least for product distributions on the hypercube, polynomials yield the best basis for L_1 regression.

Are there other suitable derandomizations of concentration inequalities? In this work, we focused on understanding the limits of k -wise independent distributions. Gopalan et al. [GKM14] gave a much more sophisticated generator with nearly optimal seed length. But could simple, natural pseudorandom distributions, such as small-bias spaces, give strong tail bounds themselves?

5 Acknowledgements

We thank Varun Kanade, Scott Linderman, Raghu Meka, Jelani Nelson, Justin Thaler, Salil Vadhan, Les Valiant, and several anonymous reviewers for helpful discussions and comments.

References

- [ABI85] Noga Alon, Laszlo Babai, and Alon Itai. A fast and simple randomized parallel algorithm for the maximal independent set problem. Technical report, Chicago, IL, USA, 1985.
- [AS04] Scott Aaronson and Yaoyun Shi. Quantum lower bounds for the collision and the element distinctness problems. *J. ACM*, 51(4):595–605, 2004.

- [Baz09] Louay M. J. Bazzi. Polylogarithmic independence can fool DNF formulas. *SIAM J. Comput.*, 38(6):2220–2272, March 2009.
- [BBC⁺01] Robert Beals, Harry Buhrman, Richard Cleve, Michele Mosca, and Ronald de Wolf. Quantum lower bounds by polynomials. *J. ACM*, 48(4):778–797, 2001.
- [Bei93] Richard Beigel. The polynomial method in circuit complexity. In *Structure in Complexity Theory Conference*, pages 82–95. IEEE Computer Society, 1993.
- [Bei94] Richard Beigel. Perceptrons, PP, and the polynomial hierarchy. *Computational Complexity*, 4:339–349, 1994.
- [Ber24] S. N. Bernstein. Le problème de l’approximation des fonctions continues sur tout l’axe réel et l’une de ses applications. *Bull. Math. Soc. France*, 52:399–410, 1924.
- [BFJ⁺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*, page 253262. ACM, 1994.
- [Bon70] Aline Bonami. Étude des coefficients de fourier des fonctions de $l^p(g)$. *Annales de l’institut Fourier*, 20(2):335–402, 1970.
- [BOW10] Eric Blais, Ryan O’Donnell, and Karl Wimmer. Polynomial regression under arbitrary product distributions. *Machine Learning*, 80(2-3):273–294, 2010.
- [BR94] M. Bellare and J. Rompel. Randomness-efficient oblivious sampling. In *FOCS*, pages 276–287, Nov 1994.
- [BRRY10] Mark Braverman, Anup Rao, Ran Raz, and Amir Yehudayoff. Pseudorandom generators for regular branching programs. *FOCS*, pages 40–47, 2010.
- [BV10] Joshua Brody and Elad Verbin. The coin problem and pseudorandomness for branching programs. In *FOCS*, pages 30–39, 2010.
- [Car51] Lennart Carleson. Bernstein’s approximation problem. *Proc. Amer. Math. Soc.*, 2:953–961, 1951.
- [Che82] E.W. Cheney. *Introduction to Approximation Theory*. AMS Chelsea Publishing Series. AMS Chelsea Pub., 1982.
- [DETT10] Anindya De, Omid Etesami, Luca Trevisan, and Madhur Tulsiani. Improved pseudorandom generators for depth 2 circuits. In Maria Serna, Ronen Shaltiel, Klaus Jansen, and Jos Rolim, editors, *Approximation, Randomization, and Combinatorial Optimization. Algorithms and Techniques*, volume 6302 of *Lecture Notes in Computer Science*, pages 504–517. 2010.
- [DFT⁺14] Dana Dachman-Soled, Vitaly Feldman, Li-Yang Tan, Andrew Wan, and Karl Wimmer. Approximate resilience, monotonicity, and the complexity of agnostic learning. *CoRR*, abs/1405.5268, 2014. To appear in SODA 2015.

- [DGJ⁺09] Ilias Diakonikolas, Parikshit Gopalan, Ragesh Jaiswal, Rocco A. Servedio, and Emanuele Viola. Bounded independence fools halfspaces. In *In Proc. 50th Annual Symposium on Foundations of Computer Science (FOCS)*, pages 171–180, 2009.
- [DLS14] Amit Daniely, Nati Linial, and Shai Shalev-Shwartz. The complexity of learning halfspaces using generalized linear methods. *CoRR*, abs/1211.0616, 2014.
- [DSTW10] Ilias Diakonikolas, Rocco A. Servedio, Li-Yang Tan, and Andrew Wan. A regularity lemma, and low-weight approximators, for low-degree polynomial threshold functions. In *Proceedings of the 2010 IEEE 25th Annual Conference on Computational Complexity, CCC '10*, pages 211–222, Washington, DC, USA, 2010. IEEE Computer Society.
- [EZ64] H. Ehlich and K. Zeller. Schwankung von polynomen zwischen gitterpunkten. *Mathematische Zeitschrift*, 86:41–44, 1964.
- [FGKP06] Vitaly Feldman, Parikshit Gopalan, Subhash Khot, and Ponnuswami. New results for learning noisy parities and halfspaces. In *Proceedings of the 47th Annual IEEE Symposium on Foundations of Computer Science, FOCS '06*, pages 563–574, Washington, DC, USA, 2006. IEEE Computer Society.
- [FK14] Vitaly Feldman and Pravesh Kothari. Agnostic learning of disjunctions on symmetric distributions. *CoRR*, abs/1405.6791, 2014.
- [FLS11] V. Feldman, H. Lee, and R. Servedio. Lower bounds and hardness amplification for learning shallow monotone formulas. *Journal of Machine Learning Research - COLT Proceedings*, 19:273292, 2011.
- [GKK08] Parikshit Gopalan, Adam Tauman Kalai, and Adam R. Klivans. Agnostically learning decision trees. In Cynthia Dwork, editor, *STOC*, pages 527–536. ACM, 2008.
- [GKM14] Parikshit Gopalan, Daniel Kane, and Raghu Meka. Pseudorandomness for concentration bounds and signed majorities. *CoRR*, abs/1411.4584, 2014.
- [GOWZ10] Parikshit Gopalan, Ryan O’Donnell, Yi Wu, and David Zuckerman. Fooling functions of halfspaces under product distributions. In *Proceedings of the 2010 IEEE 25th Annual Conference on Computational Complexity, CCC '10*, pages 223–234, Washington, DC, USA, 2010. IEEE Computer Society.
- [GR06] V. Guruswami and P. Raghavendra. Hardness of learning halfspaces with noise. In *Proceedings of FOCS 06*, page 543552, 2006.
- [HNO08] Nicholas J. A. Harvey, Jelani Nelson, and Krzysztof Onak. Sketching and streaming entropy via approximation theory. In *FOCS*, pages 489–498, 2008.
- [Hoe63] Wassily Hoeffding. Probability inequalities for sums of bounded random variables. *Journal of the American Statistical Association*, 58(301):pp. 13–30, 1963.
- [IK37] S. Izumi and T. Kawata. Quasi-analytic class and closure of $\{t^n\}$ in the interval $(-\infty, \infty)$. *Tohoku Math. J.*, 43:267–273, 1937.

- [INW94] Russell Impagliazzo, Noam Nisan, and Avi Wigderson. Pseudorandomness for network algorithms. In *STOC*, pages 356–364, 1994.
- [Kea98] M. Kearns. Efficient noise-tolerant learning from statistical queries. *Journal of the ACM*, 45(6):9831006, 1998.
- [KKM13] Daniel M. Kane, Adam Klivans, and Raghu Meka. Learning halfspaces under log-concave densities: Polynomial approximations and moment matching. In *COLT*, pages 522–545, 2013.
- [KKMS08] Adam Tauman Kalai, Adam R. Klivans, Yishay Mansour, and Rocco A. Servedio. Agnostically learning halfspaces. *SIAM J. Comput.*, 37(6):1777–1805, 2008.
- [KNP11] Michal Koucký, Prajakta Nimbhorkar, and Pavel Pudlák. Pseudorandom generators for group products. In *STOC*, pages 263–272, 2011.
- [KOS08] Adam R. Klivans, Ryan O’Donnell, and Rocco A. Servedio. Learning geometric concepts via gaussian surface area. In *FOCS*, pages 541–550, 2008.
- [KS04] Adam R. Klivans and Rocco A. Servedio. Learning DNF in time $2^{\tilde{O}(n^{1/3})}$. *J. Comput. Syst. Sci.*, 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):97114, 2007.
- [KS09] Adam R. Klivans and Alexander A. Sherstov. Cryptographic hardness for learning intersections of halfspaces. *J. Comput. Syst. Sci.*, 75(1):212, 2009.
- [KS10] Adam R. Klivans and Alexander A. Sherstov. Lower bounds for agnostic learning via approximate rank. *Computational Complexity*, 19(4):581–604, 2010.
- [KSSH94] Michael Kearns, Robert E. Schapire, Linda M. Sellie, and Lisa Hellerstein. Toward efficient agnostic learning. In *Machine Learning*, pages 341–352. ACM Press, 1994.
- [LBW95] Wee Sun Lee, Peter L. Bartlett, and Robert C. Williamson. On efficient agnostic learning of linear combinations of basis functions. In *Proceedings of the Eighth Annual Conference on Computational Learning Theory, COLT ’95*, pages 369–376, New York, NY, USA, 1995. ACM.
- [Lub07] Doron Lubinsky. A survey of weighted polynomial approximation with exponential weights. *Surveys in Approximation Theory*, 3:1–105, 2007.
- [Mar90] A. A. Markov. On a question of D. I. Mendelev. *Zapiski Imperatorskoi Akademii Nauk.*, 62:1–24, 1890.
- [MP72] Marvin Minsky and Seymour Papert. *Perceptrons: An Introduction to Computational Geometry*. MIT Press, Cambridge MA, 1972.
- [MZ10] Raghu Meka and David Zuckerman. Pseudorandom generators for polynomial threshold functions. In *Proceedings of the Forty-second ACM Symposium on Theory of Computing, STOC ’10*, pages 427–436, New York, NY, USA, 2010. ACM.

- [Nev86] Paul Nevai. Géza Freud, orthogonal polynomials and Christoffel functions. A case study. *Journal of Approximation Theory*, 48(1):3–167, 1986.
- [Nis92] Noam Nisan. $\mathcal{RL} \subset \mathcal{SC}$. In *STOC*, pages 619–623, 1992.
- [NN93] Joseph Naor and Moni Naor. Small-bias probability spaces: Efficient constructions and applications. *SIAM J. Computing*, 22:838–856, 1993.
- [NS94] N. Nisan and M. Szegedy. On the degree of boolean functions as real polynomials. *Computational Complexity*, 4:301–313, 1994.
- [NT86] Paul Nevai and Vilmos Totik. Weighted polynomial inequalities. *Constructive Approximation*, 2(1):113–127, 1986.
- [NT87] P. Nevai and V. Totik. Sharp Nikolskii inequalities with exponential weights. *Analysis Mathematica*, 13(4):261–267, 1987.
- [O’D14] Ryan O’Donnell. *Analysis of Boolean Functions*. Cambridge University Press, New York, NY, USA, 2014.
- [Pat92] Ramamohan Paturi. On the degree of polynomials that approximate symmetric boolean functions (preliminary version). In S. Rao Kosaraju, Mike Fellows, Avi Wigderson, and John A. Ellis, editors, *STOC*, pages 468–474. ACM, 1992.
- [RC66] T. J. Rivlin and E. W. Cheney. A comparison of uniform approximations on an interval and a finite subset thereof. *SIAM J. Numer. Anal.*, 3(2):311–320, 1966.
- [Rei08] Omer Reingold. Undirected connectivity in log-space. *J. ACM*, 55(4):17:1–17:24, September 2008.
- [RSV13] Omer Reingold, Thomas Steinke, and Salil Vadhan. Pseudorandomness for regular branching programs via fourier analysis. In *APPROX-RANDOM*, pages 655–670, 2013.
- [She08] Alexander A. Sherstov. Communication lower bounds using dual polynomials. *Bulletin of the EATCS*, 95:59–93, 2008.
- [She09] Alexander A. Sherstov. Separating AC^0 from depth-2 majority circuits. *SIAM J. Comput.*, 38(6):2113–2129, 2009.
- [SSS95] J. Schmidt, A. Siegel, and A. Srinivasan. Chernoff–Hoeffding bounds for applications with limited independence. *SIAM J. Discrete Mathematics*, 8(2):223–250, 1995.
- [SSSS11] S. Shalev-Shwartz, O. Shamir, and K. Sridharan. Learning kernel-based halfspaces with the 0-1 loss. *SIAM Journal on Computing*, 40:16231646, 2011.
- [SV14] Sushant Sachdeva and Nisheeth K. Vishnoi. Faster algorithms via approximation theory. *Foundations and Trends in Theoretical Computer Science*, 9(2):125–210, 2014.
- [Val84] Leslie G. Valiant. A theory of the learnable. *Commun. ACM*, 27(11):1134–1142, 1984.

[Wim10] Karl Wimmer. Agnostically learning under permutation invariant distributions. In *Proceedings of the 2010 IEEE 51st Annual Symposium on Foundations of Computer Science, FOCS '10*, pages 113–122, Washington, DC, USA, 2010. IEEE Computer Society.

A Upper Bound for Limited Independence

Theorem 5 follows from the following well-known [SSS95, BR94] lemma, which we prove for completeness.

Lemma 24. *Let $X \in \{\pm 1\}^n$ be uniform and $r \in \mathbb{R}^n$. For all even $k \geq 2$,*

$$\mathbb{E} \left[(X \cdot r)^k \right] \leq \left(e \|r\|_2^2 k \right)^{k/2}.$$

An even stronger form of Lemma 24 follows immediately from the hypercontractivity theorem [Bon70] [O'D14, §9]: Letting $f(x) = x \cdot r$, we have

$$\mathbb{E} \left[(X \cdot r)^k \right] = \|f\|_k^k \leq \left((k-1)^{\deg(f)/2} \|f\|_2 \right)^k = \left(\sqrt{k-1} \|r\|_2 \right)^k,$$

as required. A self-contained proof follows.

Proof. We start by bounding the moment generating function of $X \cdot r$: Let $t \in \mathbb{R}$ be fixed later. For any $i \in [n]$, we have

$$\mathbb{E} \left[e^{tr_i X_i} \right] = \frac{1}{2} (e^{tr_i} + e^{-tr_i}) = \sum_{k=0}^{\infty} \frac{(tr_i)^k + (-tr_i)^k}{2k!} = \sum_{k=0}^{\infty} \frac{(tr_i)^{2k}}{(2k)!} \leq \sum_{k=0}^{\infty} \frac{(t^2 r_i^2)^k}{2^k k!} = e^{t^2 r_i^2 / 2}.$$

By independence,

$$\mathbb{E} \left[e^{t(X \cdot r)} \right] = \prod_{i=1}^n \mathbb{E} \left[e^{tr_i X_i} \right] \leq \prod_{i=1}^n e^{t^2 r_i^2 / 2} = e^{t^2 \|r\|_2^2 / 2}.$$

Now we have

$$\mathbb{E} \left[e^{t(X \cdot r)} \right] = \sum_{k=0}^{\infty} \frac{t^k}{k!} \mathbb{E} \left[(X \cdot r)^k \right] \leq e^{t^2 \|r\|_2^2 / 2}.$$

We wish to bound a single moment, namely $\mathbb{E} \left[(X \cdot r)^{k_*} \right]$ for an even k_* . We do this by picking one term out of the above infinite sum. We have $\mathbb{E} \left[(X \cdot r)^k \right] \geq 0$ for even k , so these terms can be removed from the sum without increasing it. By changing the sign of t , we can ensure that the sum of the odd terms is positive and thus

$$\frac{t^{k_*}}{k_*!} \mathbb{E} \left[(X \cdot r)^{k_*} \right] \leq \sum_{k=0}^{\infty} \frac{t^k}{k!} \mathbb{E} \left[(X \cdot r)^k \right] = \mathbb{E} \left[e^{t(X \cdot r)} \right] \leq e^{t^2 \|r\|_2^2 / 2}.$$

Rearranging and setting $t = \pm \sqrt{k_*} / \|r\|_2$, we obtain

$$\mathbb{E} \left[(X \cdot r)^{k_*} \right] \leq \frac{k_*!}{t^{k_*}} e^{t^2 \|r\|_2^2 / 2} = \frac{k_*! \|r\|_2^{k_*} e^{k_*/2}}{\sqrt{k_*}^{k_*}} \leq \left(\frac{k_*^2 \|r\|_2^2 e}{k_*} \right)^{k_*/2} = (e \|r\|_2^2 k_*)^{k_*/2},$$

as required. \square

Now we can prove the upper bound for k -wise independence using the connection between moment bounds and tail bounds [SSS95].

Proof of Theorem 5. Note that, if $X \in \{\pm 1\}^n$ is k -wise independent, then

$$\mathbb{E} \left[(X \cdot r)^k \right] = \sum_{i_1 \cdots i_k \in [n]} \left(\prod_{j=1}^k r_{i_j} \right) \cdot \mathbb{E} \left[\prod_{j=1}^k X_{i_j} \right]$$

is the same as for uniform X , as this is the expectation of a degree- k polynomial. By Lemma 24 and Markov's inequality, we have (assuming k is even),

$$\mathbb{P} [|X \cdot r| \geq T] = \mathbb{P} \left[(X \cdot r)^k \geq T^k \right] \leq \frac{\mathbb{E} \left[(X \cdot r)^k \right]}{T^k} \leq \left(\frac{e \|r\|_2^2 k}{T^2} \right)^{k/2}.$$

Substituting $k = 2 \lceil \eta \log_e(1/\delta) \rceil$ and $T = e^{(\eta+1)/2\eta} \sqrt{k} \|r\|_2$, we have

$$\mathbb{P} [|X \cdot r| \geq T] \leq \left(\frac{e \|r\|_2^2 k}{(e^{(\eta+1)/2\eta} \sqrt{k} \|r\|_2)^2} \right)^{\lceil \eta \log_e(1/\delta) \rceil} = e^{-\lceil \eta \log_e(1/\delta) \rceil / \eta} \leq \delta.$$

□