

# On Agnostic PAC Learning using $\mathcal{L}_2$ -polynomial Regression and Fourier-based Algorithms

Mohsen Heidari and Wojciech Szpankowski,  
Department of Computer Science, Purdue University,  
{mheidari, szpan}@purdue.edu

## Abstract

We develop a framework using Hilbert spaces as a proxy to analyze PAC learning problems with structural properties. We consider a joint Hilbert space incorporating the relation between the true label and the predictor under a joint distribution  $D$ . We demonstrate that agnostic PAC learning with 0-1 loss is equivalent to an optimization in the Hilbert space domain. With our model, we revisit the PAC learning problem using methods based on *least-squares* such as  $\mathcal{L}_2$  polynomial regression and Linial's low-degree algorithm. We study learning with respect to several hypothesis classes such as half-spaces and polynomial-approximated classes (i.e., functions approximated by a fixed-degree polynomial). We prove that (under some distributional assumptions) such methods obtain generalization error up to  $2P_{opt}$  with  $P_{opt}$  being the optimal error of the class. Hence, we show the tightest bound on generalization error when  $P_{opt} \leq 0.2$ .

## I. INTRODUCTION

We study binary classification using polynomial regression from the agnostic PAC learning perspective [1], [2]. In this problem, multiple training instances are generated IID according to an underlying distribution  $D$  on the feature-label sets  $\mathcal{X} \times \{-1, 1\}$ . In addition, we are given a hypothesis class with respect to which the learning process takes place. If  $P_{opt}$  is the minimum error attained using the given class, then the objective of the learning algorithm is to output, with high probability, a classifier whose generalization error is not greater than  $P_{opt} + \epsilon$ .

To gain computational efficiency or analytical tractability, many conventional learning methods such as support-vector machine (SVM) rely on intermediate loss functions other than the natural 0 – 1 loss. Square loss is an example that is a basis for  $\mathcal{L}_2$ -polynomial regression or another variant of SVM known

as LS-SVM [3]. The well-known “low-degree” algorithm [4] is also known to be in this category of algorithms [5]. Such methods have been analyzed for many PAC learning problems. Under the *realizability* assumption where  $P_{opt} = 0$ , the  $\mathcal{L}_2$ -polynomial regression and the low-degree algorithm are PAC learners for a variety of hypothesis classes [6]–[8]. Under the agnostic setting where  $P_{opt} > 0$ , the current results are not that promising. The best known results for  $\mathcal{L}_2$ -polynomial regression (and the low-degree algorithm under the uniform distribution) are  $8P_{opt}$  and  $\frac{1}{4} + P_{opt}(1 - P_{opt})$  for classes such as half-spaces or polynomial-approximated classes [2], [5].

In this paper, we develop a framework using Hilbert spaces as a proxy to analyze such problems. We consider a joint Hilbert space incorporating the relation between the true label and the predictor under the joint distribution  $D$ . This is unlike conventional analysis using Hilbert spaces that focus only on the predictors with marginal  $D_x$  on the features. As a byproduct, we improve the above mentioned bounds and show that the generalization error of  $\mathcal{L}_2$ -polynomial regression and the low-degree algorithm is less than  $2P_{opt}$ . This bound improves upon the previous bounds when  $P_{opt} \leq 0.2$ . We show that methods based on square loss are suitable for learning classes with appropriate geometrical properties.

#### A. Our approach

We develop our framework by constructing two Hilbert spaces one with respect to the true underlying distribution  $D$  and the other with respect to the empirical one. The first one is  $\mathcal{L}_2(D)$ , that is all real-valued functions  $f$  on  $\mathcal{X} \times \mathcal{Y}$  such that  $\mathbb{E}[f(\mathbf{X}, Y)^2] < \infty$ . The second one is  $\mathcal{L}_2(\hat{D})$  with  $\hat{D}$  being the empirical distribution of the training set. With this formulation, the true label  $Y$  and the training labels are understood as a member of these spaces. With this formulation, the generalization error of any classifier  $c$  equals  $\frac{1}{4} \|Y - c\|_{2,D}^2$ . Similarly, when the distance is calculated in the second Hilbert space, we obtain a characterization of the empirical error. Hence, minimizing the generalization (or empirical) error is equivalent to minimizing the distance between  $Y$  and the classifier  $c$  in the first (or second) Hilbert space. We argue that the mentioned hypothesis classes have appropriate structures using that allows us to drive lower bounds on its minimum error  $P_{opt}$ . For instance, given  $k$ , the polynomial-approximated class is characterized by the subspace of  $\mathcal{L}_2(D)$  spanned by polynomials of degree up to  $k$ . With this structure, finding  $P_{opt}$  is equivalent to finding the minimum distance of  $Y$  to the subspace spanned by polynomials of degree up to  $k$ . As for the learning algorithms, we argue the low-degree algorithm and  $\mathcal{L}_2$ -polynomial regression have suitable structures using which we drive our upper bounds on their generalization errors. For instance, in the case of  $\mathcal{L}_2$  polynomial regression, the error of any classifier of the form  $\text{sign}[p(x) - \theta]$ , with  $\theta$  chosen appropriately, is bounded from above by  $\frac{1}{2} \|Y - p\|_2^2$ . Hence, minimizing the squares-loss as in  $\mathcal{L}_2$ -regression yields an error less than  $2P_{opt}$ .

## B. Summary of the Results

In this work, we first present a more general version of the low-degree algorithm incorporating non-uniform but product probability distributions. We refer to this generalization as Fourier algorithm. With our framework, we study learning with respect to three well-known hypothesis classes. The first class is half-spaces consisting of all the Boolean-valued functions of the form  $c(\mathbf{x}) = \text{sign}[\sum_{j=1}^d w_j x_j - \theta]$ . The second class is called polynomial-approximated functions. Given a positive integer  $k$  and  $\epsilon > 0$ , it consists of Boolean-valued functions that are approximated by a degree  $k$  polynomial with square error up to  $\epsilon^2$ . The thirist class is a generalization of the second. We use our framework to analyze learning these hypothesis classes using  $\mathcal{L}_2$ -polynomial regression and the Fourier algorithm. Below, is the summary of our results:

1) The  $\mathcal{L}_2$  polynomial regression with degree  $k$  outputs a hypothesis  $\hat{g}$  whose generalization error has the following properties:

- For learning polynomial-approximated classes, it is less than  $2P_{opt} + 3\epsilon$  (Theorem 1).
- For learning half-spaces, when the marginal  $D_{\mathbf{x}}$  is uniform over the unit ball in  $\mathbb{R}^d$ , it is less than  $2P_{opt} + 3\epsilon$  (Theorem 3).
- For learning *generalized concentrated classes*, under any distribution, it is less than  $2P_{opt} + \epsilon$  (Theorem 4).

2) If the marginal  $D_{\mathbf{x}}$  is a product probability distribution on  $\{-1, 1\}^d$ , then with probability  $(1 - \delta)$ , the Fourier algorithm outputs a hypothesis such that its generalization error is less than  $2P_{opt} + 2\epsilon$  for learning polynomial-approximated classes.

## C. Related Works

The low-degree algorithm is introduced by [4] with PAC learning guarantees under the uniform distribution over  $\{-1, 1\}^d$ . This algorithm which is based on the Fourier expansion on the Boolean cube has been used for in various problems [6], [8], [9]. The  $\mathcal{L}_2$  polynomial regression along with its  $\mathcal{L}_1$  counterpart is introduced by [5] for learning with respect to polynomial-approximated classes,  $k$ -juntas, and half-spaces. Learning with respect to such classes has been studied extensively in the literature [5], [10]–[12]. Among such classes, learning with respect to half-spaces is the most challenging. In the case of *proper* agnostic PAC learning, where the algorithm’s predictor must be a half-space, it is an NP-hard problem [13], [14]. Even without the *proper* restriction, the problem is NP-hard. That said, under distributional assumptions, polynomial time algorithms are introduced [5], [15], [16]. Among them are the *improper* learning algorithms based on regression methods such as  $\mathcal{L}_1$  or  $\mathcal{L}_2$  polynomial regression

[4], [5]. In particular, [5] proved that  $\mathcal{L}_1$  polynomial regression learns a range of hypothesis classes such as half-spaces (under distributional assumptions) and polynomial-approximated classes

## II. PRELIMINARIES

**Notation:** The input set is denoted by  $\mathcal{X}$  which is a subset of  $\mathbb{R}^d$  for some positive integer  $d$ . The output set is denoted by  $\mathcal{Y}$  which is a subset of  $\mathbb{R}$ . In binary classification  $\mathcal{Y} = \{-1, 1\}$ . For shorthand, the random vectors in  $\mathbb{R}^d$  are denoted by  $\mathbf{X} = (X_1, X_2, \dots, X_d)$ . Further, for any ordered subset  $\mathcal{J} = \{j_1, j_2, \dots, j_m\}$ , by  $X^{\mathcal{J}}$  denote the random vector  $(X_{j_1}, X_{j_2}, \dots, X_{j_m})$ . Similarly, by  $x^{\mathcal{J}}$  denote the vector  $(x_{j_1}, x_{j_2}, \dots, x_{j_m})$ . For a pair of functions  $f, g$  on  $\mathcal{X}$ , the notation  $f \equiv g$  means that  $f(x) = g(x)$  for all  $x \in \mathcal{X}$ . Lastly, for any natural number  $\ell$ , the set  $\{1, 2, \dots, \ell\}$  is denoted by  $[\ell]$ .

### A. A Hilbert Space Representation

We first develop a Hilbert Space formulation for the binary classification problem. Let  $D$  be a joint probability distribution on the input-output set  $\mathcal{X} \times \mathcal{Y}$ . In this paper, it is assumed that the marginal  $D_{\mathbf{x}}$  of any joint distribution  $D$  on  $\mathcal{X} \times \mathcal{Y}$  has finite moments. Consider a Hilbert space of all real-valued functions  $f : \mathcal{X} \times \mathcal{Y} \mapsto \mathbb{R}$  which are  $\mathcal{L}_2(D)$ , that is  $\mathbb{E}_D[f(\mathbf{X}, Y)^2] < \infty$ . The inner product between two members  $f, g$  is defined as

$$\langle f, g \rangle \triangleq \mathbb{E}_D[f(\mathbf{X}, Y)g(\mathbf{X}, Y)].$$

Given any integer  $p > 0$  and distribution  $D$ , the  $p$ -norm of a function  $f$  is defined as

$$\|f\|_{p,D} \triangleq (\mathbb{E}_D[f(\mathbf{X}, Y)^p])^{1/p}.$$

Given any training sample  $\mathcal{S} = \{(\mathbf{x}_i, y_i) : i = 1, 2, \dots, n\}$ , let  $\hat{D}$  denote its empirical distribution, that is a uniform distribution on  $\mathcal{S}$  and zero outside of it. Associated with this distribution, we consider the Hilbert space  $\mathcal{L}_2(\hat{D})$  with the inner product and norms defined based on the empirical distribution  $\hat{D}$ . We use this formulation to study the binary classification problem where  $\mathcal{Y} = \{-1, 1\}$ . Therefore, the generalization error of any predictor  $c : \mathcal{X} \mapsto \{-1, 1\}$  can be written in terms of the inner products as

$$\mathbb{P}_D\{Y \neq c(\mathbf{X})\} = \frac{1}{2} - \frac{1}{2}\langle Y, c \rangle_D = \frac{1}{4}\|Y - c\|_{2,D}^2, \quad (1)$$

where, with slight abuse of notation,  $Y$  is understood as the mapping  $(x, y) \mapsto y$  and  $c$  is understood as a mapping on  $\mathcal{X} \times \mathcal{Y}$  which depends only on  $\mathcal{X}$ . Similarly, the empirical error of  $c$  is equal to

$$\hat{\mathbb{P}}_{\hat{D}}\{Y \neq c(\mathbf{X})\} = \frac{1}{2} - \frac{1}{2}\langle Y, c \rangle_{\hat{D}} = \frac{1}{4}\|Y - c\|_{2,\hat{D}}^2.$$

The goal now is to derive bounds on the minimum generalization error when learning with respect to various hypothesis classes. In Section III we describe  $\mathcal{L}_2$ -polynomial regression and the Fourier algorithm,

in Section IV we study polynomial-approximated classes, and finally in Section V we discuss half-spaces, and more general hypothesis classes that have structural properties.

### III. PAC LEARNING WITH $\mathcal{L}_2$ -POLYNOMIAL REGRESSION

We employ a PAC learning algorithm using  $\mathcal{L}_2$ -polynomial regression. Given a training set, the objective of the polynomial regression is to minimize the empirical square loss over all polynomials of degree up to  $k$ . This process can be implemented by stochastic gradient descent or by solving a linear system of equations. We describe how this polynomial regression can be used for PAC learning. Let  $\hat{p}$  be the output of the polynomial regression. The idea is to shift the polynomial  $\hat{p}$  by a threshold  $\theta$  and take its sign. This process is demonstrated as Algorithm 1.

---

**Algorithm 1** PAC Learning with  $\mathcal{L}_2$ -Polynomial Regression

---

**Input:** Degree parameter  $k$ , and training samples  $\{(\mathbf{x}(i), y(i)), i \in [n]\}$ .

1: Find a polynomial  $\hat{p}$  of degree up to  $k$  that minimizes

$$\frac{1}{n} \sum_i (y(i) - p(\mathbf{x}(i)))^2.$$

2: Find  $\theta \in [-1, 1]$  such that the empirical error of  $\text{sign}[\hat{p}(\mathbf{x}) - \theta]$  is minimized.

3: **return**  $\hat{g} \equiv \text{sign}[\hat{p} - \theta]$ .

---

#### A. Fourier-Based Learning Algorithm

We present another variant of  $\mathcal{L}_2$  polynomial regression, known as the low-degree (Fourier) algorithm [4]. Although this algorithm is more efficient than the polynomial regression, it requires binary input set  $\mathcal{X} = \{-1, 1\}^d$ . The low-degree algorithm was originally designed for uniform distribution on the Boolean cube. In this paper, we present a more general version of it for incorporating non-uniform but product probability distributions on  $\{-1, 1\}^d$  [17]. In this approach, the objective is to find an estimate of the  $p^*$  polynomial that minimizes the square loss  $\|Y - p^*\|_{2,D}$  under the true distribution. This method is based on the Fourier expansion on the Boolean cube [18] and is summarized in the following.

Under product probability distribution on  $\{-1, 1\}^d$ , any bounded real-valued functions can be written as

$$f(\mathbf{x}) = \sum_{S \subseteq [d]} f_S \psi_S(\mathbf{x}),$$

where  $f_{\mathcal{S}}$ 's are the Fourier coefficients and calculated as  $f_{\mathcal{S}} \triangleq \langle f, \psi_{\mathcal{S}} \rangle$  for every subset  $\mathcal{S} \subseteq [d]$ . Further, the parity  $\psi_{\mathcal{S}}$  is a monomial defined as

$$\psi_{\mathcal{S}}(x) = \prod_{j \in \mathcal{S}} \frac{x_j - \mu_j}{\sigma_j},$$

with  $\mu_j$  and  $\sigma_j$  being the mean and standard deviation of the  $X_j$ , respectively. As the distribution is unknown, these quantities are estimated in the algorithm.

As a result, we can write the Fourier decomposition of the optimal polynomial  $p^*$ . For that, we have the following statement:

**Fact 1.** *Let  $D$  be a probability distribution with the marginal  $D_{\mathbf{x}}$  that is a product probability distribution on  $\{-1, 1\}^d$ . Then, the optimal polynomial  $p^*$  admits the following Fourier decomposition*

$$p^* \equiv \sum_{\mathcal{S} \subseteq [d]: |\mathcal{S}| \leq k} \langle Y, \psi_{\mathcal{S}} \rangle \psi_{\mathcal{S}}.$$

With that decomposition, the idea behind the Fourier algorithm is to compute an empirical estimate of  $\langle Y, \psi_{\mathcal{S}} \rangle$ . This is demonstrated as Algorithm 2.

---

**Algorithm 2** Fourier-Based Learning

---

**Input:** Training samples  $\{(\mathbf{x}(i), y(i)), i \in [n]\}$ .

- 1: Compute the empirical mean  $\hat{\mu}_j$  and standard deviation  $\hat{\sigma}_j$  of each feature.
- 2: For every  $\mathcal{S} \subseteq [d]$  with  $|\mathcal{S}| \leq k$ , construct the empirical parity as  $\hat{\psi}_{\mathcal{S}}(\mathbf{x}) = \prod_{j \in \mathcal{S}} \frac{x_j - \hat{\mu}_j}{\hat{\sigma}_j}$ .
- 3: Compute the empirical Fourier coefficients  $a_{\mathcal{S}}$ , for every  $\mathcal{S}$  with at most  $k$  elements, as

$$a_{\mathcal{S}} = \frac{1}{n} \sum_{i=1}^n y(i) \hat{\psi}_{\mathcal{S}}(\mathbf{x}(i)).$$

- 4: Construct and return the function  $\hat{\Pi}_Y$  as

$$\hat{\Pi}_Y(\mathbf{x}) \triangleq \sum_{\mathcal{S}: |\mathcal{S}| \leq k} a_{\mathcal{S}} \hat{\psi}_{\mathcal{S}}(\mathbf{x}).$$


---

In the following lemma which is proved in Appendix A, we derive bounds for estimating the optimal polynomial  $p^*$ .

**Lemma 1.** *Let  $D$  be a probability distribution with the marginal  $D_{\mathbf{x}}$  that is a product probability distribution on  $\{-1, 1\}^d$ . Given  $\delta \in (0, 1)$ , with probability at least  $(1 - \delta)$ , the following inequality holds*

$$\|p^* - \hat{\Pi}_Y\|_2 \leq O\left(\sqrt{\frac{d^k c_k}{(k-1)!n} \log \frac{4d^k}{(k-1)!\delta}}\right), \quad (2)$$

where  $c_k \triangleq \max_{\mathcal{S} \subseteq [d], |\mathcal{S}| \leq k} \|\psi_{\mathcal{S}}\|_{\infty}^2$  and  $n$  is the number of samples.

#### IV. POLYNOMIALLY APPROXIMATED CLASS

In this section, we study agnostic PAC learning with respect to concept classes whose members are approximated by fixed-degree polynomials. We adopt the Hilbert space representation in Section II-A to analyze PAC learning using Algorithm 1 and 2. We start with the following formulation:

**Definition 1.** Given  $\epsilon \in [0, 1]$ ,  $k \in \mathbb{N}$  and any probability distribution  $D_{\mathbf{X}}$  on  $\mathcal{X}$ , a concept class  $\mathcal{C}$  of functions  $c : \mathcal{X} \mapsto \{-1, 1\}$  is  $(\epsilon, k)$ -approximated if

$$\sup_{c \in \mathcal{C}} \inf_{p \in \mathcal{P}_k} \mathbb{E}[(c(\mathbf{X}) - p(\mathbf{X}))^2] \leq \epsilon^2,$$

where  $\mathcal{P}_k$  is the set of all polynomials of degree up to  $k$ .

We consider agnostic PAC learning with respect to  $\mathcal{C}$  and under the 0 – 1 loss function. The minimum generalization error and empirical error of  $\mathcal{C}$  are, respectively, defined as

$$\begin{aligned} P_{opt} &\triangleq \min_{c \in \mathcal{C}} \mathbb{P}_D \{Y \neq c(\mathbf{X})\}, \\ \hat{P}_{opt} &\triangleq \min_{c \in \mathcal{C}} \hat{\mathbb{P}}_{\hat{D}} \{Y \neq c(\mathbf{X})\}. \end{aligned}$$

We use the Hilbert space representation in Section II-A and provide a lower bound on  $P_{opt}$ .

**Lemma 2.** *The minimum generalization error attainable by any  $(\epsilon, k)$  concept class  $\mathcal{C}$  is bounded from below as*

$$P_{opt} \geq \frac{1}{2} - \frac{1}{2} \|p^*\|_{1,D} - \epsilon,$$

where  $p^* = \arg \min_{p \in \mathcal{P}_k} \mathbb{E}_D[(Y - p(\mathbf{X}))^2]$ .

*Proof.* From (1) the 0 – 1 loss of any function  $c \in \mathcal{C}$  can be written as  $\mathbb{P}\{Y \neq c(\mathbf{X})\} = \frac{1}{2} - \frac{1}{2} \langle Y, c \rangle$ . Let  $p \in \mathcal{P}_k$  be such that  $\|c - p\|_{2,D} \leq \epsilon$ . Then, by adding and subtracting  $p$ , we obtain that

$$\begin{aligned} \langle Y, c \rangle &= \langle Y, p \rangle + \langle Y, (c - p) \rangle \\ &\leq \langle Y, p \rangle + \|Y\|_2 \|c - p\|_2 \leq \langle Y, p \rangle + \epsilon, \end{aligned} \tag{3}$$

where the first inequality follows from Cauchy–Schwarz inequality and the second inequality follows because  $\|Y\|_2 = 1$ . Note that  $\mathcal{P}_k$ , the set of all polynomials on  $\mathcal{X}$  with degree upto  $k$ , is a (finite dimensional) subspace inside the Hilbert space  $\mathcal{L}_2(D)$ . Therefore, it has an orthonormal basis denoted by  $\{\Psi_1, \Psi_2, \dots, \Psi_m\}$ , where  $m$  is less than  $O(d^k)$ . As a result, the polynomial  $p$  can be written as  $p \equiv \sum_{j=1}^m \langle p, \Psi_j \rangle \Psi_j$ . Hence,

$$\langle Y, p \rangle = \sum_{j=1}^m \langle p, \Psi_j \rangle \langle Y, \Psi_j \rangle = \langle \Pi_Y, p \rangle,$$

where  $\Pi_Y \equiv \sum_{j=1}^m \langle Y, \Psi_j \rangle \Psi_j$  is the *projection* of  $Y$  onto this subspace. Consequently, from the above equality and (3), we obtain that

$$\begin{aligned} \langle Y, c \rangle &\leq \langle \Pi_Y, p \rangle + \epsilon = \langle \Pi_Y, c \rangle + \langle \Pi_Y, (p - c) \rangle + \epsilon \\ &\leq \langle \Pi_Y, c \rangle + \|\Pi_Y\|_2 \|p - c\|_2 + \epsilon \\ &\leq \|\Pi_Y\|_1 + \|\Pi_Y\|_2 \|p - c\|_2 + \epsilon \\ &\leq \|\Pi_Y\|_1 + 2\epsilon, \end{aligned}$$

where the second inequality follows from Cauchy–Schwarz inequality, the third one holds as  $|c(\mathbf{x})| = 1$  and the last inequality follows from Bessel’s inequality, implying  $\|\Pi_Y\|_2 \leq 1$ , and the assumption that  $\|p - c\|_2 \leq \epsilon$ . Next, we proceed with the following fact about the projection.

**Fact 2.**  $\Pi_Y$  the projection of  $Y$  onto  $\mathcal{P}_k$  is the polynomial minimizing  $\mathbb{E}[(Y - p(\mathbf{X}))^2]$  over all  $p \in \mathcal{P}_k$ .

The proof is complete by the following fact implying that  $\Pi_Y \equiv p^*$ . □

We show in Section III-A that the lower-bound in Lemma 2 helps to prove our results for the low-degree algorithm.

#### A. PAC Learning Bounds

Next, we analyze Algorithm 1 and 2 for this class and prove the first main result of the paper.

**Theorem 1.** *Given  $\epsilon > 0$  and  $k \in \mathbb{N}$ , the degree  $k$   $\mathcal{L}_2$  polynomial regression agnostically PAC learns any  $(\epsilon, k)$ -approximated concept class with expected error up to*

$$2P_{\text{opt}} + 3\epsilon + \sqrt{\frac{2 d^{k+1}}{n} \log \frac{en}{d^{k+1}}},$$

where  $d$  is the number of input variables and  $n$  is the sample size.

*Proof.* To derive an upper bound on the empirical error of  $\hat{g}$ , we first consider a weaker version of the algorithm. The idea is to select  $\theta$  randomly instead of optimizing it as in the algorithm. For that, we establish the following lemma.

**Lemma 3.** *Suppose  $\theta$  is a random variable with the probability density function  $f_\theta(t) = 1 - |t|$ , for  $t \in [-1, 1]$ . Then, the following bound holds for any polynomial  $p$*

$$\mathbb{E}_\theta \left[ \hat{\mathbb{P}} \left\{ Y \neq \text{sign}[p(\mathbf{X}) - \theta] \right\} \right] \leq \frac{1}{2} \|Y - p\|_{2, \hat{D}}^2.$$



*Proof.* Note that  $y \neq \text{sign}(p(\mathbf{x}) - \theta)$ , if  $\theta$  is between  $y$  and  $p(\mathbf{x})$ . Hence, the expected empirical error of  $\text{sign}[p(\mathbf{X}) - \theta]$  with respect to the random  $\theta$  equals to

$$\begin{aligned} & \mathbb{E}_\theta \left[ \hat{\mathbb{P}} \left\{ Y \neq \text{sign}[p(\mathbf{X}) - \theta] \right\} \right] \\ &= \frac{1}{n} \sum_i \mathbb{E}_\theta \left[ \mathbb{1} \{ y_i \neq \text{sign}(p(\mathbf{x}_i) - \theta) \} \right] \\ &= \frac{1}{n} \sum_i \underbrace{\mathbb{P} \left\{ \theta \in [p(\mathbf{x}_i), y_i] \cup [y_i, p(\mathbf{x}_i)] \right\}}_{\mathbb{P}_i}. \end{aligned} \quad (4)$$

Next, we show that  $\mathbb{P}_i \leq \frac{1}{2}(y_i - p(\mathbf{x}_i))^2$  for all  $(\mathbf{x}_i, y_i)$ 's. Suppose  $y_i = 1$ . If  $p(\mathbf{x}_i) > 1$ , then  $\mathbb{P}_i = 0$  as  $\theta \leq 1$ . If  $p(\mathbf{x}_i) \in [0, 1]$ , then

$$\begin{aligned} \mathbb{P}_i &= \mathbb{P} \left\{ \theta \in [p(\mathbf{x}_i), 1] \right\} = \int_{p(\mathbf{x}_i)}^1 (1 - t) dt \\ &= \frac{1}{2} (1 - p(\mathbf{x}_i))^2 = \frac{1}{2} (y_i - p(\mathbf{x}_i))^2. \end{aligned}$$

If  $p(\mathbf{x}_i) \in [-1, 0]$ , then

$$\begin{aligned} \mathbb{P}_i &= \mathbb{P} \left\{ \theta \in [p(\mathbf{x}_i), 1] \right\} = \int_{p(\mathbf{x}_i)}^1 1 - |t| dt \\ &= \frac{1}{2} + \int_{p(\mathbf{x}_i)}^0 (1 + t) dt \\ &= \frac{1}{2} - p(\mathbf{x}_i) - \frac{1}{2} (p(\mathbf{x}_i))^2 \\ &\leq \frac{1}{2} (1 + |p(\mathbf{x}_i)|)^2 = \frac{1}{2} (y_i - p(\mathbf{x}_i))^2. \end{aligned}$$

Lastly, if  $p(\mathbf{x}_i) < -1$ , then  $\mathbb{P}_i = 1$  because  $\theta \geq -1$ . In this case also  $\mathbb{P}_i \leq \frac{1}{2}(y_i - p(\mathbf{x}_i))^2$ . The case for  $y_i = -1$  follows by symmetricity. Hence, we obtain the following inequality

$$\mathbb{E}_\theta \left[ \hat{\mathbb{P}} \left\{ Y \neq \hat{g}(\mathbf{X}) \right\} \right] \leq \frac{1}{n} \sum_i \frac{1}{2} (y_i - p(\mathbf{x}_i))^2.$$

The proof is complete by noting that the right-hand side equals to  $\frac{1}{2} \|Y - p\|_{2, \hat{D}}^2$ .  $\square$

Consequently, from the lemma and due the fact that  $\theta$  in the algorithm is selected to minimize the empirical error, we obtain that

$$\hat{\mathbb{P}} \left\{ Y \neq \hat{g}(\mathbf{X}) \right\} \leq \frac{1}{2} \|Y - \hat{p}\|_{2, \hat{D}}^2, \quad (5)$$

where  $\hat{p}$  is the output of  $\mathcal{L}_2$ -polynomial regression and  $\hat{g} \equiv \text{sign}[\hat{p} - \theta]$ , as in Algorithm 1. Let  $c^*$  be the predictor with minimum generalization error in the  $(\epsilon, k)$ -approximated concept class. Let  $p$  be a degree

$k$  polynomial such that  $\|c^* - p\|_2 \leq \epsilon$ . Since  $\hat{p}$  minimizes the empirical 2-norm, then the right-hand side of (5) satisfies

$$\frac{1}{2}\|Y - \hat{p}\|_{2,\hat{D}}^2 \leq \frac{1}{2}\|Y - p^*\|_{2,\hat{D}}^2. \quad (6)$$

We proceed by taking the expected error of the empirical error with respect to the random training samples. From (5) and (6) we obtain the following inequalities

$$\begin{aligned} \mathbb{E}\left[\hat{\mathbb{P}}\{Y \neq \hat{g}(\mathbf{X})\}\right] &\leq \frac{1}{2}\mathbb{E}\left[\|Y - p^*\|_{2,\hat{D}}^2\right] = \frac{1}{2}\|Y - p^*\|_{2,D}^2 \\ &\stackrel{(a)}{\leq} \frac{1}{2}\left(\|Y - c^*\|_{2,D} + \|p^* - c^*\|_{2,D}\right)^2 \\ &\leq \frac{1}{2}\left(\|Y - c^*\|_{2,D} + \epsilon\right)^2 \\ &\stackrel{(b)}{\leq} \frac{1}{2}\left(\|Y - c^*\|_{2,D}^2 + 4\epsilon + \epsilon^2\right) \\ &\stackrel{(c)}{\leq} 2\mathbf{P}_{opt} + \frac{5}{2}\epsilon, \end{aligned} \quad (7)$$

where (a) holds from Minkowski's inequality for 2-norm, (b) holds as  $\|Y - c^*\|_{2,D} \leq 2$ , and (c) holds because of the second equality in (1) and that  $\mathbf{P}_{opt} = \mathbb{P}\{Y \neq c^*(\mathbf{X})\}$ .

Next, we connect the empirical error of  $\hat{g}$  to its generalization error. Note that the Vapnik–Chervonenkis (VC) dimension of all functions of the form  $\text{sign}[p]$  for some polynomial of degree upto  $k$  does not exceed  $d^{k+1}$ . Therefore, from VC theory ( See Corollary 3.19 in [19]) for any  $\delta$ , with probability at least  $(1 - \delta)$ , the following inequality holds

$$\begin{aligned} \mathbb{P}\{Y \neq \hat{g}(\mathbf{X})\} &\leq \hat{\mathbb{P}}\{Y \neq \hat{g}(\mathbf{X})\} + \sqrt{\frac{2 d^{k+1}}{n} \log \frac{en}{d^{k+1}}} \\ &\quad + \sqrt{\frac{\log \frac{1}{\delta}}{2n}}. \end{aligned} \quad (8)$$

Set  $\delta = \exp\{-\frac{1}{2}n\epsilon^2\}$ . Therefore, the proof is complete by taking the expectation and combining it with the last bound in (7).  $\square$

**PAC bounds for the Fourier algorithm:** Next, we employ a low-degree (Fourier) algorithm (Algorithm 2) for PAC learning with respect to the polynomially approximated hypothesis class.

**Theorem 2.** *Let  $D$  be a joint probability distribution with marginal  $D_X$  that is a product probability distribution on  $\{-1, 1\}^d$ . Then, for any  $\delta \in [0, 1]$ , with probability at least  $1 - \delta$ , the Fourier-based algorithm agnostically PAC learns any  $(\epsilon, k)$ -approximated concept class with generalization error up to*

$$2\mathbf{P}_{opt} + 2\epsilon + O\left(\sqrt{\frac{d^k c_k}{(k-1)!n} \log \frac{4d^k}{(k-1)!\delta}}\right), \quad (9)$$

where  $c_k \triangleq \max_{S \subseteq [d], |S| \leq k} \|\psi_S\|_\infty^2$ .

*Proof.* We prove the theorem by characterizing the effect of 2-norm estimation on the error probability.

Let

$$p^* = \arg \min_{p \in \mathcal{P}_k} \|Y - p\|_{2,D}^2.$$

From the second equality in (1), the generalization error of  $\hat{g}$  in Algorithm 2 satisfies

$$\begin{aligned} \mathbb{P}\{Y \neq \hat{g}(\mathbf{X})\} &= \frac{1}{4} \|Y - \hat{g}\|_{2,D}^2 \\ &\leq \frac{1}{4} \left( \|Y - p^*\|_{2,D} + \|p^* - \hat{g}\|_{2,D} \right)^2 \\ &\leq \frac{1}{2} \left( \|Y - p^*\|_{2,D}^2 + \|p^* - \hat{g}\|_{2,D}^2 \right), \end{aligned} \quad (10)$$

where the inequality follows from Minkowski's inequality for 2-norm. Observe that (10) is an upper bound on the generalization error in terms of 2-norm quantities. Since  $p^*$  minimizes the square loss, the first term in (10) equals

$$\|Y - p^*\|_{2,D}^2 = 1 - \|p^*\|_2^2.$$

We proceed by bounding the second term in (10). From Minkowski's inequality for 2-norm and by adding and subtracting  $\hat{\Pi}_Y$  as in Algorithm 2, we have that

$$\begin{aligned} \|p^* - \hat{g}\|_2^2 &\leq \|p^* - \hat{\Pi}_Y\|_2^2 + \|\hat{\Pi}_Y - \hat{g}\|_2^2 \\ &\quad + 2\|p^* - \hat{\Pi}_Y\|_2 \|\hat{\Pi}_Y - \hat{g}\|_2. \end{aligned} \quad (11)$$

The first term in (11) is bounded from Lemma 1. As a result,  $\|p^* - \hat{\Pi}_Y\|_2 = \epsilon'_n$ , where

$$\epsilon'_n = O\left(\sqrt{\frac{d^k c_k}{(k-1)!n} \log \frac{4d^k}{(k-1)!\delta}}\right),$$

with probability at least  $(1-\delta)$ . As for the second term in (11), we use the identity  $|h - \text{sign}[h]| = |1 - |h||$  for any function  $h$ . Therefore, as  $\hat{g} = \text{sign}[\hat{\Pi}_Y]$ , we obtain that

$$\begin{aligned} \|\hat{\Pi}_Y - \hat{g}\|_2^2 &= \mathbb{E}\left[(1 - |\hat{\Pi}_Y(\mathbf{X})|)^2\right] \\ &= 1 + \|\hat{\Pi}_Y\|_2^2 - 2\|\hat{\Pi}_Y\|_1. \end{aligned} \quad (12)$$

Next, we show that the third term in (11) is less than  $4\epsilon'_n$ . It suffices to show that  $\|\hat{\Pi}_Y - \hat{g}\|_2 \leq 2$ . For that, we use the equality in (12). By removing the last term in (12) and taking the square root we have

$$\|\hat{\Pi}_Y - \hat{g}\|_2 \leq \sqrt{1 + \|\hat{\Pi}_Y\|_2^2}.$$

From the Minkowski's inequality we have that

$$\begin{aligned}\|\hat{\Pi}_Y\|_2 &\leq \|p^*\|_2 + \|p^* - \hat{\Pi}_Y\|_2 \\ &\leq \|p^*\|_2 + \epsilon'_n \leq 1 + \epsilon'_n.\end{aligned}$$

Hence, we get the desired bound  $\|\hat{\Pi}_Y - \hat{g}\|_2 \leq \sqrt{1 + (1 + \epsilon'_n)^2} \leq 2$ , assuming that  $\epsilon'_n \leq 1/3$ . Combining the bounds for each term in (11), we get

$$\begin{aligned}\|p^* - \hat{g}\|_2^2 &\leq \epsilon'_n{}^2 + 1 + \|\hat{\Pi}_Y\|_2^2 - 2\|\hat{\Pi}_Y\|_1 + 4\epsilon'_n \\ &\leq 1 + \|\hat{\Pi}_Y\|_2^2 - 2\|\hat{\Pi}_Y\|_1 + 5\epsilon'_n.\end{aligned}$$

We plug this inequality in (10). After rearranging the terms by adding and subtracting  $\|p^*\|_1$ , we obtain the following inequality

$$\begin{aligned}\mathbb{P}\{Y \neq \hat{g}(\mathbf{X})\} &\leq \frac{1}{2} \left( 2 - 2\|p^*\|_1 + 5\epsilon'_n \right. \\ &\quad \left. + 2(\|p^*\|_1 - \|\hat{\Pi}_Y\|_1) + (\|\hat{\Pi}_Y\|_2^2 - \|p^*\|_2^2) \right) \\ &\leq 2P_{opt} + 2\epsilon + 5\epsilon'_n,\end{aligned}$$

where the last inequality follows from Lemma 2 and the following argument for bounding the last two terms in the first inequality:

For the 1-norm difference, the Minkowski's inequality for 1-norm gives

$$\|p^*\|_1 - \|\hat{\Pi}_Y\|_1 \leq \|p^* - \hat{\Pi}_Y\|_1 \leq \|p^* - \hat{\Pi}_Y\|_2 = \epsilon'_n,$$

where the last inequality follows from the Jensen's inequality implying that  $\|\cdot\|_1 \leq \|\cdot\|_2$ .

For the difference of square of 2-norms, we apply the Minkowski's inequality for 2-norm and obtain

$$\begin{aligned}\|\hat{\Pi}_Y\|_2^2 &\leq \|p^*\|_2^2 + \|p^* - \hat{\Pi}_Y\|_2^2 + 2\|p^*\|_2\|p^* - \hat{\Pi}_Y\|_2 \\ &\leq \|p^*\|_2^2 + 3\epsilon'_n.\end{aligned}$$

where the last inequality holds as  $\|p^*\|_2 \leq 1$ . □

We end this section by presenting a simplified result of Theorem 2.

**Corollary 1.** *If the expected value of each  $X_j$  satisfies  $|\mu_j| \leq 1 - \frac{1}{k}$ , then the generalization error of the Fourier algorithm is upper bounded by*

$$2P_{opt} + 2\epsilon + O\left(\sqrt{\frac{\sqrt{k}(ed)^k}{n} \left(k \log \frac{ed}{k} + \log \frac{2\sqrt{k}}{\delta}\right)}\right).$$

*Proof.* From the definition of  $c_k$ , we can write

$$\begin{aligned} c_k &= \max_{\mathcal{S}:|\mathcal{S}|\leq k} \max_{\mathbf{x}\in\{-1,1\}^d} |\psi_{\mathcal{S}}(\mathbf{x})|^2 \leq \max_{\mathcal{S}:|\mathcal{S}|\leq k} \prod_{j\in\mathcal{S}} \frac{(1+|\mu_j|)^2}{\sigma_j^2} \\ &= \max_{\mathcal{S}:|\mathcal{S}|\leq k} \prod_{j\in\mathcal{S}} \frac{1}{1-|\mu_j|}, \end{aligned}$$

where the first inequality is from the definition of  $\psi_{\mathcal{S}}$ . The last quality holds as  $\sigma_j^2 = 1 - \mu_j^2$ . Therefore, under the assumption that  $|\mu_j| \leq 1 - \frac{1}{k}$ , the following inequality holds

$$c_k \leq \max_{\mathcal{S}:|\mathcal{S}|\leq k} \prod_{j\in\mathcal{S}} k \leq k^k$$

Hence,

$$\frac{c_k}{(k-1)!} \leq \frac{k k^k}{k!} = 2^{(k+1)\log_2 k - \log_2 k!}$$

From Stirling approximation  $\log_2 k! = k \log_2 k - k \log_2 e + O(\log_2 k)$ . Hence,

$$\frac{c_k}{(k-1)!} \leq 2^{k \log_2 e + O(\log_2 k)} = e^k + O(k).$$

Using the above inequality and Theorem 2 in the main text, we obtain the corollary. □

## V. LEARNING OTHER HYPOTHESIS CLASSES

In this section, we extend our results to two other type of concept classes. The first one is called *half-spaces* and the other one is a generalized version of the concentrated hypothesis classes.

### A. Half-spaces

In this section, we consider learning another class of functions called half-spaces. More precisely, a half-space a Boolean-valued function of the form

$$c(\mathbf{x}) = \text{sign}\left[a_0 + \sum_{j=1}^d a_j x_j\right], \quad \forall \mathbf{x} \in \mathbb{R}^d$$

where  $a_j \in \mathbb{R}$ . We start with a lower-bound on the optimal classification error of the class.

**Lemma 4.** *Let  $D$  be any joint probability distribution on  $\mathbb{R}^d \times \{-1, 1\}$  with marginal  $D_{\mathbf{x}}$  that is the uniform distribution on  $\mathbb{S}^{d-1}$  or jointly Gaussian on  $\mathbb{R}^d$ . Then, for any  $\epsilon > 0$ , the minimum generalization error of learning with respect to half-spaces satisfy the following lower bound*

$$P_{opt} \geq \frac{1}{2} - \frac{1}{2} \|p_{\epsilon}^*\|_{1, D_{\mathbf{x}}} - \epsilon,$$

where  $p_{\epsilon}^*$  is a polynomial of degree up to  $O(\frac{1}{\epsilon})$  minimizing  $\|Y - p\|_{2, D}$  among all such polynomials.

The proof of the lemma follows from Lemma 2 and [5]’s result (Theorem 6) on the sign function. This result is stated as

**Lemma 5** ([5]). *Let  $X$  be a random variable with uniform distribution on  $\mathbb{S}^{d-1}$  or jointly Gaussian on  $\mathbb{R}^d$ . Then, for any  $\epsilon > 0$ , there exists a polynomial  $p$  of degree  $O(\frac{1}{\epsilon^4})$  such that  $\mathbb{E}\left[\left(p(\mathbf{X}) - \text{sign}(\mathbf{X})\right)^2\right] \leq \epsilon^2$ .*

This lemma makes a connection between half-spaces and the polynomial-approximated class. That said, in the following theorems we show our results for PAC learning using Algorithm 1.

**Theorem 3.** *Let  $D$  be any joint probability distribution on  $\mathbb{R}^d \times \{-1, 1\}$ , with marginal  $D_X$  that is uniform on the unit sphere or jointly Gaussian. Then,  $\mathcal{L}_2$ -polynomial regression PAC learns half-spaces with expected generalization error up to*

$$2P_{opt} + 3\epsilon + \sqrt{\frac{d^{O(\frac{1}{\epsilon^4})}}{n} \log \frac{n}{d^{O(\frac{1}{\epsilon^4})}}.$$

### B. Generalized approximated class

Lastly, we finish this paper by extending our results to a more general hypothesis class. Fix a set of functions  $e_1(\mathbf{x}), e_2(\mathbf{x}), \dots, e_m(\mathbf{x})$  and let  $\mathcal{H}$  be a Hilbert space spanned by these functions. Let  $\mathcal{C}$  be a class of functions each of which is approximated by elements of  $\mathcal{H}$  with square error up to  $\epsilon$ , that is,

$$\inf_{h \in \mathcal{H}} \|c - h\|_{2,D} \leq \epsilon,$$

for any  $c \in \mathcal{C}$ . As a special case, suppose  $e_i$ ’s are all the functions of the form  $e(\mathbf{x}) = \prod_{j \in [d]} x_j^{\alpha_j}$  where  $\alpha_j$ ’s are non-negative integers adding up to  $k$ . Then  $\mathcal{C}$  is a  $(k, \epsilon)$ -approximated class as in Section IV.

**Theorem 4.** *Suppose  $A$  is any algorithm that given  $n$  training instances finds a function  $\hat{h} \in \mathcal{H}$  so that the empirical loss  $\|Y - \hat{h}\|_{2,\hat{D}}$  is minimized. Then, the predictor  $\text{sign}[\hat{h}]$  learns  $\mathcal{C}$  with expected generalization error up to*

$$2P_{opt} + 3\epsilon + O\left(\sqrt{\frac{\text{VC}(\mathcal{C})}{n} \log \frac{n}{\text{VC}(\mathcal{C})}}\right),$$

where  $\text{VC}(\mathcal{C})$  is the VC dimension of  $\mathcal{C}$ .

### ACKNOWLEDGEMENT

This work was supported in part by NSF Center on Science of Information Grants CCF-0939370 and NSF Grants CCF-1524312, CCF-2006440, CCF-2007238, and Google Research Award.

APPENDIX A  
PROOF OF LEMMA 1

a) *Mean and variance estimations*:: We first take into account the effect of the imperfections in mean and variance estimation. For tractability of our analysis, we use a fraction of the training samples just for the mean and variance estimations. As a measure of accuracy of the estimations, we require the differences  $|\hat{\mu}_j - \mu_j|$  and  $|1 - \frac{\sigma}{\hat{\sigma}}|$  to be sufficiently small with probability close to one. This is a deviation from standard measures of estimations in which the variance of the differences are required to be small. In the following lemma, we bound the estimation errors in terms of the number of the samples.

**Lemma 6.** *Given  $\epsilon_0, \delta_0 \in (0, 1)$  the following inequalities hold with probability at least  $(1 - \delta_0)$*

$$|\hat{\mu}_j - \mu_j| \leq \epsilon_0, \quad \left|1 - \frac{\sigma_j}{\hat{\sigma}_j}\right| \leq \frac{2\epsilon_0}{\sigma_j^2}, \quad (13)$$

for all  $j \in [d]$ , provided that at least  $n_0(\epsilon_0, \delta_0) = \frac{2}{\epsilon_0^2} \log \frac{2d}{\delta_0}$  samples are available.

*Proof.* Form McDiarmid's inequality, for each  $j \in [d]$  we have

$$\mathbb{P}\{|\hat{\mu}_j - \mu_j| \geq \epsilon_0\} \leq 2 \exp\left\{-\frac{n\epsilon_0^2}{2}\right\}.$$

Therefore, applying the union bound gives

$$\mathbb{P}\left\{\bigcup_{j=1}^d \{|\hat{\mu}_j - \mu_j| \geq \epsilon_0\}\right\} \leq 2d \exp\left\{-\frac{n\epsilon_0^2}{2}\right\}.$$

Thus, the right-hand side of the above inequality is less than  $\delta_0$ , if  $n \geq \frac{2}{\epsilon_0^2} \log(\frac{2d}{\delta_0})$ . As a result we obtain the inequalities for the estimation of  $\mu_j$ 's. Next, we prove the inequalities for the estimation of  $\sigma_j$ 's. For any fixed  $\hat{\mu} \in (-1, 1)$ , define the function  $h_{\hat{\mu}}(x) = \frac{\sqrt{1-x^2}}{\sqrt{1-\hat{\mu}^2}}$ . From Taylor's theorem, there exists  $\zeta \in (-1, 1)$  which is between  $x$  and  $\hat{\mu}$  such that

$$h_{\hat{\mu}}(x) = 1 - \frac{\zeta(x - \hat{\mu})}{\sqrt{(1 - \zeta^2)(1 - \hat{\mu}^2)}}.$$

As a result,

$$|h_{\hat{\mu}}(x) - 1| = \frac{|\zeta||x - \hat{\mu}|}{\sqrt{(1 - \zeta^2)(1 - \hat{\mu}^2)}} \leq \frac{|x - \hat{\mu}|}{\sqrt{(1 - (\max\{x, \hat{\mu}\})^2)(1 - \hat{\mu}^2)}}.$$

Now by setting  $x = \mu_j$  and that  $|\hat{\mu}_j - \mu_j| \leq \epsilon_0$ , we have

$$\left|\frac{\sigma_j}{\hat{\sigma}_j} - 1\right| = |h_{\hat{\mu}}(\mu_j) - 1| \leq \frac{\epsilon_0}{\hat{\sigma} \min\{\hat{\sigma}, \sigma\}}.$$

Note that,  $|\hat{\mu}_j| \leq |\mu_j| + \epsilon_0$ . Therefore,

$$\hat{\sigma}_j^2 \geq 1 - (|\mu_j| + \epsilon_0)^2 \geq \sigma_j^2 - 2\epsilon_0|\mu_j| - \epsilon_0^2 \geq \sigma_j^2 - 3\epsilon_0.$$

As a result,

$$\left| \frac{\sigma_j}{\hat{\sigma}_j} - 1 \right| \leq \frac{\epsilon_0}{\sigma_j^2 - 3\epsilon_0} \leq \frac{2\epsilon_0}{\sigma_j^2},$$

which completes the proof of the lemma.  $\square$

Now we proceed with the proof of the lemma. Let  $\bar{\Pi}_Y$  denote the version of  $\hat{\Pi}_Y$  under the assumption that  $\hat{\mu}_j = \mu_j$  and  $\hat{\sigma}_j = \sigma_j$  for all  $j \in [d]$ . Also, let  $B$  be the event that the inequalities in (13) hold. From Minkowsky's inequality, by adding and subtracting  $\bar{\Pi}_Y$  we have

$$\|\Pi_Y - \hat{\Pi}_Y\|_2 \leq \underbrace{\|\Pi_Y - \bar{\Pi}_Y\|_2}_V + \underbrace{\|\bar{\Pi}_Y - \hat{\Pi}_Y\|_2}_W.$$

Let  $V$  and  $W$  denote the first and the second term above, respectively. We proceed by the following lemmas.

**Lemma 7.** *Given any  $\delta > 0$ , the inequality  $\|\Pi_Y - \bar{\Pi}_Y\|_2 \leq \sqrt{\frac{2d^k c_k}{(k-1)!n} \log \frac{2d^k}{(k-1)! \delta}}$  holds with probability  $(1 - \delta)$ .*

*Proof.* Recall that  $\bar{\Pi}_Y$  is defined as

$$\bar{\Pi}_Y(x^d) \triangleq \sum_{\mathcal{S}:|\mathcal{S}|\leq k} \bar{f}_{\mathcal{S}} \psi_{\mathcal{S}}(x^d),$$

where the Fourier-estimates  $\bar{f}_{\mathcal{S}}$  are defined as  $\bar{f}_{\mathcal{S}} \triangleq \frac{1}{n} \sum_i Y(i) \psi_{\mathcal{S}}(X(i))$ . In addition, by definition of the projection function  $\Pi_Y$ , we have

$$\Pi_Y(\mathbf{x}) = \sum_{\mathcal{S}:|\mathcal{S}|\leq k} f_{\mathcal{S}} \psi_{\mathcal{S}}(\mathbf{x}), \quad \forall \mathbf{x} \in \mathcal{X}^d.$$

Therefore, from Parseval's identity, the 2-norm factors as

$$\|\Pi_Y - \bar{\Pi}_Y\|_2^2 = \sum_{\mathcal{S}:|\mathcal{S}|\leq k} |f_{\mathcal{S}} - \bar{f}_{\mathcal{S}}|^2.$$

In what follows, we show that  $|f_{\mathcal{S}} - \bar{f}_{\mathcal{S}}| \leq \epsilon$  for all subsets  $\mathcal{S} \subseteq [d]$  with  $|\mathcal{S}| \leq k$ . Note that  $\bar{f}_{\mathcal{S}}$  is a function of the training random samples  $(X(i), Y(i)), i = 1, 2, \dots, n$ . Observe that  $\mathbb{E}[\bar{f}_{\mathcal{S}}] = f_{\mathcal{S}}$  which implies that  $\bar{f}_{\mathcal{S}}$  is an unbiased estimation of  $f_{\mathcal{S}}$ . Since the samples are drawn independent and identically distributed (i.i.d.), we apply McDiarmid's inequality to bound the probability of the event  $|f_{\mathcal{S}} - \bar{f}_{\mathcal{S}}| \geq \epsilon'$ .



For that, fix  $i \in [d]$  and suppose  $(\mathbf{X}(i), Y(i))$  in the training set is replaced with an i.i.d. copy  $(\tilde{\mathbf{X}}(i), \tilde{Y}(i))$ . With this replacement  $\bar{f}_S$  is changed to another random variable denoted by  $\tilde{f}_S$ . Then

$$\begin{aligned} |\bar{f}_S - \tilde{f}_S| &= \frac{1}{n} |Y(i)\psi_S(\mathbf{X}(i)) - \tilde{Y}(i)\psi_S(\tilde{\mathbf{X}}(i))| \\ &\leq \frac{1}{n} |Y(i)\psi_S(\mathbf{X}(i))| + |\tilde{Y}(i)\psi_S(\tilde{\mathbf{X}}(i))| \\ &\leq \frac{1}{n} |\psi_S(\mathbf{X}(i))| + |\psi_S(\tilde{\mathbf{X}}(i))| \\ &\leq \frac{2}{n} \|\psi_S\|_\infty, \end{aligned}$$

where  $\|\psi_S\|_\infty = \max_{\mathbf{x}} |\psi_S(\mathbf{x})|$ . Let  $c_k = \max_{S \subseteq [d], |S| \leq k} \|\psi_S\|_\infty^2$ . Then, from McDiarmid's inequality, for any  $\epsilon' \in (0, 1)$

$$\mathbb{P}\left\{ \max_{S: |S| \leq k} |\bar{f}_S - f_S| \geq \epsilon' \right\} \leq 2 \left[ \sum_{m=0}^k \binom{d}{m} \right] \exp\left\{ -\frac{n\epsilon'^2}{2c_k} \right\}, \quad (14)$$

where we also used the union bound. For  $k \leq d/2$ , we obtain that

$$\sum_{m=0}^k \binom{d}{m} \leq k \frac{d^k}{k!}.$$

As a result, with probability at least  $(1 - \delta)$ ,  $\max_{S: |S| \leq k} |\bar{f}_S - f_S| \leq \sqrt{\frac{2c_k}{n} \log \frac{2d^k}{(k-1)!\delta}}$ . Hence, we with probability at least  $(1 - \delta)$

$$\|\Pi_Y - \bar{\Pi}_Y\|_2^2 \leq \frac{2d^k c_k}{(k-1)!n} \log \frac{2d^k}{(k-1)!\delta},$$

and the proof is complete by taking the square root of both sides.  $\square$

**Lemma 8.** *Conditioned on  $B$ , the inequalities  $\|\bar{\Pi}_Y - \hat{\Pi}_Y\|_\infty \leq \lambda(\epsilon)$  hold, almost surely, for all  $k$ -element subsets  $\mathcal{J} \subset [d]$ , where  $\lambda$  is a function satisfying  $\lambda(\epsilon_0) = O\left(\frac{kd^k c_k}{(k-1)!} \epsilon_0\right)$  as  $\epsilon_0 \rightarrow 0$ .*

Recall that the function  $\bar{\Pi}_Y$  is defined as

$$\bar{\Pi}_Y(x^d) \triangleq \sum_{S: |S| \leq k} \bar{f}_S \psi_S(x^d),$$

where the Fourier-estimates  $\bar{f}_S$  are defined as

$$\bar{f}_S \triangleq \frac{1}{n} \sum_i Y(i) \psi_S(X(i)).$$

From triangle inequality for  $\infty$ -norm and the definition of  $\hat{\Pi}_Y$  and  $\bar{\Pi}_Y$  we obtain

$$\|\hat{\Pi}_Y - \bar{\Pi}_Y\|_\infty \leq \sum_{S: |S| \leq k} \|\hat{f}_S \hat{\psi}_S - \bar{f}_S \psi_S\|_\infty. \quad (15)$$

Again by triangle inequality and by adding and subtracting  $\bar{f}_S \widehat{\psi}_S$ , we obtain that

$$\begin{aligned} \|\widehat{f}_S \widehat{\psi}_S - \bar{f}_S \psi_S\|_\infty &\leq \|\widehat{f}_S \widehat{\psi}_S - \bar{f}_S \widehat{\psi}_S\|_\infty + \|\bar{f}_S \widehat{\psi}_S - \bar{f}_S \psi_S\|_\infty \\ &= |\widehat{f}_S - \bar{f}_S| \|\widehat{\psi}_S\|_\infty + |\bar{f}_S| \|\widehat{\psi}_S - \psi_S\|_\infty. \end{aligned}$$

Next, note that from triangle inequality

$$|\widehat{f}_S - \bar{f}_S| \leq \frac{1}{n} \sum_i |\widehat{\psi}_S(\mathbf{x}(i)) - \psi_S(\mathbf{x}(i))| \leq \|\psi_S - \widehat{\psi}_S\|_\infty.$$

Therefore,

$$\|\widehat{f}_S \widehat{\psi}_S - \bar{f}_S \psi_S\|_\infty \leq (\|\widehat{\psi}_S\|_\infty + |\bar{f}_S|) \|\widehat{\psi}_S - \psi_S\|_\infty. \quad (16)$$

We proceed by bounding each term above. As for the first term we have, that  $\|\widehat{\psi}_S\|_\infty \leq \|\psi_S\|_\infty + \|\widehat{\psi}_S - \psi_S\|_\infty$ . As for the second term, we have

$$\bar{f}_S = \frac{1}{n} \sum_i Y(i) \psi_S(\mathbf{X}(i)) \leq \|\psi_S\|_\infty.$$

Lastly, the third term is bounded using the following lemma.

**Lemma 9.** *Conditioned on  $B$ , the inequality  $\|\psi_S - \widehat{\psi}_S\|_\infty \leq \gamma(\epsilon_0)$  holds, almost surely, where  $\gamma$  is a function satisfying  $\gamma(\epsilon_0) = O(k\epsilon_0\sqrt{c_k})$  as  $\epsilon_0 \rightarrow 0$ .*

Before proving this lemma, we complete our argument. As a result of this lemma and using the triangle inequality, we obtain from (16) that

$$\begin{aligned} \|\widehat{f}_S \widehat{\psi}_S - \bar{f}_S \psi_S\|_\infty &\leq (2\|\psi_S\|_\infty + \|\widehat{\psi}_S - \psi_S\|_\infty) \|\widehat{\psi}_S - \psi_S\|_\infty \\ &\leq (2\sqrt{c_k} + \gamma(\epsilon_0))\gamma(\epsilon_0). \end{aligned}$$

Lastly, from (15) we get the following bound

$$\|\widehat{\Pi}_Y - \bar{\Pi}_Y\|_\infty \leq \lambda(\epsilon_0) \triangleq \frac{d^k}{(k-1)!} (2\sqrt{c_k}\gamma(\epsilon_0) + \gamma^2(\epsilon_0)).$$

It is not difficult to check that  $\lambda(\epsilon_0) = O(\frac{kd^k c_k}{(k-1)!} \epsilon_0)$  as  $\epsilon_0 \rightarrow 0$ . Now it remains to prove Lemma 9 which is given below:

**Proof of Lemma 9:** We start with the triangle inequality for  $\infty$ -norm by adding and subtracting  $b_S \psi_S$ :

$$\|\psi_S - \widehat{\psi}_S\|_\infty \leq \|\psi_S - b_S \psi_S\|_\infty + \|b_S \psi_S - \widehat{\psi}_S\|_\infty.$$

Note that  $b_S \psi_S \equiv \prod_{j \in S} \frac{x_j - \mu_j}{\hat{\sigma}_j}$ . Now, using the triangle inequality on the second term above, we have

$$\begin{aligned}
\|b_S \psi_S - \hat{\psi}_S\|_\infty &= \|b_S \psi_S \pm \left( \sum_{l \in S} \prod_{j < l} \frac{x_j - \hat{\mu}_j}{\hat{\sigma}_j} \prod_{r > l} \frac{x_r - \mu_r}{\hat{\sigma}_r} \right) - \hat{\psi}_S\|_\infty \\
&\leq \sum_{l \in S} \frac{|\mu_l - \hat{\mu}_l|}{\hat{\sigma}_l} \left\| \prod_{j < l} \frac{(x_j - \hat{\mu}_j)}{\hat{\sigma}_j} \prod_{r > l} \frac{(x_r - \mu_r)}{\hat{\sigma}_r} \right\|_\infty \\
&\leq \frac{\epsilon}{\sigma_{\min}} \sum_{l \in S} \left\| \prod_{j < l} \frac{(x_j - \hat{\mu}_j)}{\hat{\sigma}_j} \prod_{r > l} \frac{(x_r - \mu_r)}{\hat{\sigma}_r} \right\|_\infty \\
&\leq \frac{\epsilon}{\sigma_{\min}} \sum_{l \in S} \prod_{j < l} \frac{(1 + |\hat{\mu}_j|)}{\hat{\sigma}_j} \prod_{r > l} \frac{(1 + |\mu_r|)}{\hat{\sigma}_r} \\
&\stackrel{(a)}{\leq} \frac{\epsilon}{\sigma_{\min}} \sum_{l \in S} \prod_{j < l} \frac{(1 + |\mu_j|)(1 + \epsilon)}{\hat{\sigma}_j} \prod_{r > l} \frac{(1 + |\mu_r|)}{\hat{\sigma}_r} \\
&\stackrel{(b)}{\leq} \frac{\epsilon}{\sigma_{\min}} b_S \sum_{l \in S} \prod_{j \in S} \frac{(1 + |\mu_j|)(1 + \epsilon)}{\sigma_j} \\
&\stackrel{(c)}{\leq} \frac{k\epsilon}{\sigma_{\min}} b_S (1 + \epsilon)^k \|\psi_S\|_\infty,
\end{aligned}$$

where (a) follows from the inequality  $(1 + |\hat{\mu}_j|) \leq (1 + |\mu_j|)(1 + \epsilon)$ , and (b) follows from  $(1 + |\mu_j|) \leq (1 + |\mu_j|)(1 + \epsilon)$ . Lastly, (c) holds as  $|\mathcal{S}| \leq k$  and because  $\|\psi_S\|_\infty = \prod_{j \in S} \frac{1 + |\mu_j|}{\sigma_j}$ .

$$\|\psi_S - \hat{\psi}_S\|_\infty \leq |1 - b_S| \|\psi_S\|_\infty + \frac{k\epsilon}{\sigma_{\min}} b_S (1 + \epsilon)^k \|\psi_S\|_\infty. \quad (17)$$

From the assumption of the lemma and the definition of  $b_S$  we obtain that

$$1 - (1 + \epsilon)^{|\mathcal{S}|} \leq 1 - b_S \leq 1 - (1 - \epsilon)^{|\mathcal{S}|}.$$

Since  $\epsilon \in (0, 1)$  and  $|\mathcal{S}| \leq k$ , then  $(1 - \epsilon)^{|\mathcal{S}|} \geq 1 - k\epsilon$ . Also, from the fact that  $(1 + x) \leq e^x$  for all  $x \in \mathbb{R}$ , we obtain

$$1 - e^{k\epsilon} \leq 1 - b_S \leq k\epsilon \leq e^{k\epsilon} - 1. \quad (18)$$

Lastly, combining (17) and (18) gives the following inequality

$$\|\psi_S - \hat{\psi}_S\|_\infty \leq (e^{k\epsilon} - 1) \|\psi_S\|_\infty + \frac{k\epsilon}{\sigma_{\min}} (1 + \epsilon)^{2k} \|\psi_S\|_\infty.$$

The proof is complete by noting that  $\|\psi_S\|_\infty \leq \sqrt{c_k}$ . ■

From Lemma 8, we know that  $W$  is measurable with respect to  $B$ . In particular, conditioned on  $B$ ,  $W \leq \lambda(\epsilon_0)$ . Therefore, from the above lemmas and using the inequality  $\|\cdot\|_2 \leq \|\cdot\|_\infty$ , we have, with probability  $(1 - \delta_0)(1 - \delta)$  that

$$\|\Pi_Y - \bar{\Pi}_Y\|_2 \leq \sqrt{\frac{2d^k c_k}{(k-1)!n} \log \frac{2d^k}{(k-1)! \delta}} + \lambda(\epsilon_0).$$

Now set  $\epsilon_0 = \sqrt{\frac{2}{n_0} \log \frac{2d}{\delta}}$  with  $\delta_0 = \delta$ . Then, with  $n_0 = O(n)$ , we get with probability  $(1 - \delta)^2$  that  $\|\Pi_Y - \bar{\Pi}_Y\|_2 = O\left(\sqrt{\frac{2d^k c_k}{(k-1)!n} \log \frac{2d^k}{(k-1)!\delta}}\right)$ . Now the proof is complete by changing  $\delta$  to  $\delta/2$  and noting that  $(1 - \delta/2)^2 \geq 1 - \delta$ .  $\square$

## REFERENCES

- [1] L. G. Valiant, “A theory of the learnable,” *Communications of the ACM*, vol. 27, no. 11, pp. 1134–1142, nov 1984.
- [2] M. J. Kearns, R. E. Schapire, and L. M. Sellie, “Toward efficient agnostic learning,” *Machine Learning*, vol. 17, no. 2-3, pp. 115–141, 1994.
- [3] J. Suykens and J. Vandewalle, *Neural Processing Letters*, vol. 9, no. 3, pp. 293–300, 1999.
- [4] N. Linial, Y. Mansour, and N. Nisan, “Constant depth circuits, Fourier transform, and learnability,” *J. ACM*, vol. 40, no. 3, pp. 607–620, 1993.
- [5] A. T. Kalai, A. R. Klivans, Yishay Mansour, and R. A. Servedio, “Agnostically learning halfspaces,” in *Proc. 46th Annual IEEE Symp. Foundations of Computer Science (FOCS’05)*, Oct. 2005, pp. 11–20.
- [6] E. Mossel, R. O’Donnell, and R. A. Servedio, “Learning functions of  $k$  relevant variables,” *J. Comput. Syst. Sci.*, vol. 69, no. 3, pp. 421–434, 2004.
- [7] E. Mossel, R. O’Donnell, and R. P. Servedio, “Learning juntas,” in *Proc. ACM Symp. on Theory of Computing*, 2003, pp. 206–212.
- [8] E. Blais, R. O’Donnell, and K. Wimmer, “Polynomial regression under arbitrary product distributions,” *Machine learning*, vol. 80, no. 2-3, pp. 273–294, 2010.
- [9] M. Heidari, G. I. Shamir, and W. Szpankowski, “Fourier-based universal learning,” *Journal of Machine Learning Research (JMLR)*, 2020.
- [10] A. R. Klivans, P. M. Long, and R. A. Servedio, “Learning halfspaces with malicious noise,” *Journal of Machine Learning Research*, vol. 10, no. 12, 2009.
- [11] A. Birnbaum and S. S. Shwartz, “Learning halfspaces with the zero-one loss: time-accuracy tradeoffs,” in *Advances in Neural Information Processing Systems*, 2012, pp. 926–934.
- [12] I. Diakonikolas, T. Gouleakis, and C. Tzamos, “Distribution-independent pac learning of halfspaces with massart noise,” in *Advances in Neural Information Processing Systems*, 2019, pp. 4749–4760.
- [13] A. Klivans and P. Kothari, “Embedding hard learning problems into gaussian space,” in *Approximation, Randomization, and Combinatorial Optimization. Algorithms and Techniques (APPROX/RANDOM 2014)*. Schloss Dagstuhl-Leibniz-Zentrum fuer Informatik, 2014.
- [14] V. Guruswami and P. Raghavendra, “Hardness of learning halfspaces with noise,” *SIAM Journal on Computing*, vol. 39, no. 2, pp. 742–765, 2009.
- [15] P. Awasthi, M. F. Balcan, and P. M. Long, “The power of localization for efficiently learning linear separators with noise,” *Journal of the ACM*, vol. 63, no. 6, pp. 1–27, feb 2017.
- [16] A. Daniely, “A ptas for agnostically learning halfspaces,” in *Conference on Learning Theory*, 2015, pp. 484–502.
- [17] M. L. Furst, J. C. Jackson, and S. W. Smith, “Improved learning of  $AC^0$  functions,” in *COLT*, vol. 91, 1991, pp. 317–325.
- [18] R. O’Donnell, *Analysis of boolean functions*. Cambridge University Press, 2014.
- [19] M. N. Y. U. Mohri, A. (Google, I. Rostamizadeh, A. U. of California, and B. Talwalkar, *Foundations of Machine Learning*. MIT Press Ltd, 2018.