

---

# Attribute-Efficient PAC Learning of Low-Degree Polynomial Threshold Functions with Nasty Noise

---

Shiwei Zeng<sup>1</sup> Jie Shen<sup>1</sup>

## Abstract

The concept class of low-degree polynomial threshold functions (PTFs) plays a fundamental role in machine learning. In this paper, we study PAC learning of  $K$ -sparse degree- $d$  PTFs on  $\mathbb{R}^n$ , where any such concept depends only on  $K$  out of  $n$  attributes of the input. Our main contribution is a new algorithm that runs in time  $(nd/\epsilon)^{O(d)}$  and under the Gaussian marginal distribution, PAC learns the class up to error rate  $\epsilon$  with  $O(\frac{K^{4d}}{\epsilon^{2d}} \cdot \log^{5d} n)$  samples even when an  $\eta \leq O(\epsilon^d)$  fraction of them are corrupted by the nasty noise of Bshouty et al. (2002), possibly the strongest corruption model. Prior to this work, attribute-efficient robust algorithms are established only for the special case of sparse homogeneous halfspaces. Our key ingredients are: 1) a structural result that translates the attribute sparsity to a sparsity pattern of the Chow vector under the basis of Hermite polynomials, and 2) a novel attribute-efficient robust Chow vector estimation algorithm which uses exclusively a restricted Frobenius norm to either certify a good approximation or to validate a sparsity-induced degree- $2d$  polynomial as a filter to detect corrupted samples.

## 1. Introduction

A polynomial threshold function (PTF)  $f : \mathbb{R}^n \rightarrow \{-1, 1\}$  is of the form  $f(x) = \text{sign}(p(x))$  for some  $n$ -variate polynomial  $p$ . The class of low-degree PTFs plays a fundamental role in learning theory owing to its remarkable power for rich representations (Mansour, 1994; Anthony & Bartlett, 1999; Hellerstein & Servedio, 2007; O'Donnell, 2014). In this paper, we study *attribute-efficient* learning of degree-

<sup>1</sup>Department of Computer Science, Stevens Institute of Technology, Hoboken, New Jersey, USA. Correspondence to: Shiwei Zeng <szeng4@stevens.edu>, Jie Shen <jie.shen@stevens.edu>.

*Proceedings of the 40<sup>th</sup> International Conference on Machine Learning*, Honolulu, Hawaii, USA. PMLR 202, 2023. Copyright 2023 by the author(s).

$d$  PTFs: if the underlying PTF  $f^*$  is promised to depend only on at most  $K$  unknown attributes of the input, whether and how can one learn  $f^*$  by collecting  $O(\text{poly}(K^d, \log n))$  samples under the classic probably approximately correct (PAC) learning model (Valiant, 1984).

A very special case of the problem, the class of sparse halfspaces (i.e. degree-1 PTFs), has been extensively investigated in machine learning and statistics (Littlestone, 1987; Blum, 1990; Gentile, 2003; Plan & Vershynin, 2013a), and a fruitful set of results have been established even under strong noise models (Plan & Vershynin, 2013b; Awasthi et al., 2016; Zhang, 2018; Zhang et al., 2020; Shen & Zhang, 2021). It, however, turns out that the theoretical and algorithmic understanding of learning sparse degree- $d$  PTFs fall far behind the linear counterpart.

In the absence of noise, the problem can be cast as solving a linear program by thinking of degree- $d$  PTFs as halfspaces on the space expanded by all monomials of degree at most  $d$  (Maass & Turán, 1994). However, the problem becomes subtle when samples might be contaminated adversarially. In this work, we consider the nasty noise of Bshouty et al. (2002), perhaps the strongest noise model in classification.

**Definition 1** (PAC learning with nasty noise). Denote by  $\mathcal{H}_{d,K}$  the class of  $K$ -sparse degree- $d$  PTFs on  $\mathbb{R}^n$ . Let  $D$  be a distribution on  $\mathbb{R}^n$  and  $f^* \in \mathcal{H}_{d,K}$  be the underlying PTF. A nasty adversary  $\text{EX}(\eta)$  takes as input a sample size  $N$  requested by the learner, draws  $N$  instances independently according to  $D$  and annotates them by  $f^*$ , to form a clean sample set  $\bar{S} = \{(x_i, f^*(x_i))\}_{i=1}^N$ . The adversary may then inspect the learning algorithm and uses its unbounded computational power to replace at most an  $\eta$  fraction with carefully constructed samples for some  $\eta < \frac{1}{2}$ , and returns the corrupted set  $\tilde{S}'$  to the learner. The goal of the learner is to output a concept  $\hat{f} : \mathbb{R}^n \rightarrow \{-1, 1\}$ , such that with probability  $1 - \delta$  (over the randomness of the samples and all internal random bits of the learning algorithm),  $\Pr_{x \sim D}(\hat{f}(x) \neq f^*(x)) \leq \epsilon$  for any prescribed error rate  $\epsilon \in (0, 1)$  and failure probability  $\delta \in (0, 1)$ . We say an algorithm PAC learns the class  $\mathcal{H}_{d,K}$  if the guarantee holds uniformly for any member  $f^* \in \mathcal{H}_{d,K}$ .

Bshouty et al. (2002) presented a computationally inefficient

algorithm for learning any concept class with near-optimal sample complexity and noise tolerance as far as the concept class has finite VC-dimension (see Theorem 7 therein). Since the VC-dimension of  $\mathcal{H}_{d,K}$  is  $O(K^d \log n)$ , we have:

**Fact 2.** There exists an inefficient algorithm that PAC learns  $\mathcal{H}_{d,K}$  with near-optimal sample complexity  $O(K^d \log n)$  and noise tolerance  $\Omega(\epsilon)$ .

Designing efficient algorithms that match such statistical guarantees thus becomes a core research line.

On the one hand, for the special case of homogeneous sparse halfspaces, [Shen & Zhang \(2021\)](#) gave a state-of-the-art algorithm with sample complexity  $O(K^2 \log^5 n)$  and noise tolerance  $\Omega(\epsilon)$  when the instance distribution  $D$  is isotropic log-concave. On the other hand, for learning general sparse low-degree PTFs, very little is known, since the structure of PTFs is tremendously complex. To our knowledge, it appears that the only known approach is to reducing the problem to a generic approach proposed in the early work of [Kearns & Li \(1988\)](#). In particular, Theorem 12 therein implies that any concept class  $\mathcal{H}$  can be PAC learned with nasty noise in polynomial time provided that there exists a polynomial-time algorithm that PAC learns it in the absence of noise and that  $\eta \leq O\left(\frac{\epsilon}{\text{VCdim}(\mathcal{H})} \log \frac{\text{VCdim}(\mathcal{H})}{\epsilon}\right)$ , where  $\text{VCdim}(\mathcal{H})$  denotes the VC-dimension of  $\mathcal{H}$ . We therefore have the following (see Appendix B for the proof):

**Fact 3.** There exists an efficient algorithm that draws  $C_0 \cdot \frac{K^{3d} \log^3 n}{\epsilon^6} \log \frac{1}{\delta}$  samples from the nasty adversary and PAC learns  $\mathcal{H}_{d,K}$  provided that  $\eta \leq O\left(\frac{d\epsilon}{K^d} \log \frac{1}{\epsilon}\right)$ , where  $C_0 > 0$  is an absolute constant.

The result above is appealing since it makes no distributional assumption and it runs in polynomial time. However, the main issue is on the noise tolerance: the robustness of the algorithm degrades significantly when  $K$  is large. For example, in the interesting regime  $K = \Theta(\log n)$ , the noise tolerance is dimension-dependent, meaning that the algorithmic guarantees are brittle in high-dimensional problems. See [Kalai et al. \(2005\)](#); [Klivans et al. \(2009\)](#); [Long & Servedio \(2011\)](#); [Awasthi et al. \(2017\)](#); [Shen \(2021b\)](#); [Shen & Zhang \(2021\)](#) and a comprehensive survey by [Diakonikolas & Kane \(2019\)](#) for the importance and challenges of obtaining dimension-independent noise tolerance.

## 1.1. Main results

Throughout the paper, we always assume:

**Assumption 1.**  $D$  is the Gaussian distribution  $\mathcal{N}(0, I_{n \times n})$ .

Our main result is an attribute-efficient algorithm that runs in time  $(nd/\epsilon)^{O(d)}$  and PAC learns  $\mathcal{H}_{d,K}$  with *dimension-independent* noise tolerance.

**Theorem 4** (Theorem 19, informal). *Assume that  $D$  is the standard Gaussian distribution  $\mathcal{N}(0, I_{n \times n})$ . There is an*

*algorithm that runs in time  $(nd/\epsilon)^{O(d)}$  and PAC learns  $\mathcal{H}_{d,K}$  by drawing  $C \cdot \frac{K^{4d} (d \log n)^{5d}}{\epsilon^{2d+2}}$  samples from the nasty adversary for some absolute constant  $C > 0$ , provided that  $\eta \leq O(\epsilon^{d+1}/d^{2d})$ .*

**Remark 5** (Sample complexity). It is known that for *efficient and outlier-robust* algorithms,  $\Omega(K^2)$  samples are necessary to obtain an error bound of  $O(\epsilon)$  even for linear models [\(Diakonikolas et al., 2017\)](#). Thus, the multiplicative factor  $K^{4d}$  in our sample complexity bound is very close to the best possible scaling of  $K^{2d}$  and the best known result in Fact 3. The exponent  $d$  in the factor  $\frac{1}{\epsilon^{2d+2}}$  comes from our two-step approach: we will first robustly estimate the Chow vector [\(Chow, 1961\)](#) of  $f^*$  up to error  $\epsilon_0$  using  $\Omega(1/\epsilon_0^2)$  samples, and then apply an algorithmic result of [Trevisan et al. \(2009\)](#); [De et al. \(2014\)](#); [Diakonikolas et al. \(2018a\)](#) to construct a PTF with misclassification error rate of  $O(d \cdot \epsilon_0^{1/d+1})$ . This in turn suggests that we have to set  $\epsilon_0 = (\epsilon/d)^{d+1}$  in order to get the target error rate  $\epsilon$ . As noted in [Diakonikolas et al. \(2018a\)](#), such overhead on the scaling of  $\epsilon$  is inherent when using only Chow vector to establish PAC guarantees for degree- $d$  PTFs.

**Remark 6** (Noise tolerance). Our noise tolerance matches the best known one given by [Diakonikolas et al. \(2018a\)](#), which studied non-sparse low-degree PTFs. When the degree  $d$  is a constant, the noise tolerance reads as  $\epsilon^{\Omega(1)}$ , qualitatively matching the information-theoretic limit of  $\Omega(\epsilon)$ . Yet, the existence of efficient low-degree PTF learners with optimal noise tolerance is widely open.

**Remark 7** (Comparison to prior works). Most in line with this work are [Shen & Zhang \(2021\)](#) and [Diakonikolas et al. \(2018a\)](#). [Shen & Zhang \(2021\)](#) gave a state-of-the-art algorithm for learning  $K$ -sparse halfspaces, but their algorithm cannot be generalized to learn sparse low-degree PTFs. [Diakonikolas et al. \(2018a\)](#) developed an efficient algorithm for learning non-sparse PTFs. Their sample complexity bound reads as  $\frac{1}{\epsilon^{2d+2}} \cdot (nd)^{O(d)}$ , which is not attribute-efficient and thus is inapplicable for real-world problems where the number of samples is orders of magnitude less than that of attributes. At a high level, our result can be thought of as a significant generalization of both.

**Remark 8** (Running time). The computational cost of our algorithm is  $(nd/\epsilon)^{O(d)}$  which we believe may not be significantly improved in the presence of the nasty noise. This is because the adversary has the power to inspect the algorithm and to corrupt any samples, which forces any robust algorithm to carefully verify the covariance matrix whose size is  $n^{O(d)} \times n^{O(d)}$ ; see Section 1.2 and Section 3. Under quite different problem settings, prior works leveraged the underlying sparsity for improved computational complexity; see Section 1.3.

## 1.2. Overview of main techniques

Our starting point to learn sparse low-degree PTFs is an elegant algorithmic result from Diakonikolas et al. (2018a), which shows that as far as one is able to approximate the Chow vector of  $f^*$  (Chow, 1961), it is possible to reconstruct a PTF  $\hat{f}$  with PAC guarantee in time  $(nd/\epsilon)^{O(d)}$ . To apply such scheme, we will need to 1) properly define the Chow vector since it depends on the choice of the basis of polynomials; and 2) estimate the Chow vector of  $f^*$  in time  $(nd/\epsilon)^{O(d)}$ . Our technical novelty lies in new structural and algorithmic ingredients to achieve attribute efficiency.

### 1) Structural result: attribute sparsity induces sparse Chow vectors under the basis of Hermite polynomials.

Prior works such as the closely related work of Diakonikolas et al. (2018a) tend to use the basis of monomials to define Chow vectors. However, there is no guarantee that such definition would exhibit the desired sparsity structure. For example, consider that  $D$  is the standard Gaussian distribution. For  $K$ -sparse degree- $d$  PTFs on  $\mathbb{R}^n$ , the number of monomials with non-zero coefficients can be as large as  $(\frac{n}{d})^{\lfloor \frac{d}{2} \rfloor}$  for any  $K \leq n/2$ ,  $d \geq 2$  (see Lemma 22). Our first technical insight is that, the condition that a degree- $d$  PTF is  $K$ -sparse implies a  $k$ -sparse Chow vector with respect to the basis of Hermite polynomials of degree at most  $d$  ( $k$  is roughly  $2K^d$ , see Lemma 10). This endows a sparsity structure of the Chow vector of  $f^*$ , which in turn is leveraged into our algorithmic design since we can now focus on a much narrower space of Chow vectors and thus lower sample complexity.

It is worth mentioning that while we present our results under Gaussian distribution and thus use the Hermite polynomials as the basis to ease analysis, the choice of basis can go well beyond that; see Appendix B.4 for more discussions.

### 2) Algorithmic result: attribute-efficient robust Chow vector estimation.

Denote by  $m(x)$  the vector of all  $n$ -variate Hermite polynomials of degree at most  $d$ . The Chow vector (also known as Fourier coefficients) of a Boolean-valued function  $f : \mathbb{R}^n \rightarrow \{-1, 1\}$  is defined as  $\chi_f := \mathbb{E}_{x \sim D}[f(x) \cdot m(x)]$ . As we discussed,  $\chi_{f^*}$  is  $k$ -sparse on an unknown support set. To estimate it within error  $\epsilon_0$  in  $\ell_2$ -norm, it suffices to find some  $k$ -sparse vector  $u$ , such that for all  $2k$ -sparse unit vector  $v$ , we have  $|\langle v, u - \chi_{f^*} \rangle| \leq \epsilon_0$  (see Lemma 45). We will choose  $u$  as an empirical Chow vector, i.e.  $u = \sum_{(x,y) \in \bar{S}''} y \cdot m(x)$ , where  $\bar{S}''$  needs to be a carefully selected subset of  $\bar{S}'$ . Now recent developments in noise-tolerant classification (Awasthi et al., 2017; Shen & Zhang, 2021; Shen, 2023) suggest that such estimation error is governed by the maximum eigenvalue on all possible  $2k$ -sparse directions of the empirical covariance matrix<sup>1</sup>  $\Sigma :=$

<sup>1</sup>We will slightly abuse the terminology of covariance matrix to refer to the one without subtracting the mean.

$\frac{1}{|\bar{S}''|} \sum_{x \in \bar{S}''} m(x)m(x)^\top$ . This structural result can be cast into algorithmic design: find a large enough subset  $\bar{S}''$  such that the maximum sparse eigenvalue of  $(\Sigma - I)$  is close to zero (note that  $\mathbb{E}_{x \sim D}[m(x)m(x)^\top] = I$ ). Unfortunately, there are two technical challenges: first, computing the maximum sparse eigenvalue is NP-hard; second, searching for such a subset is also computationally intractable.

**2a) Small Frobenius norm certifies a good approximation.** We tackle the first challenge by considering a sufficient condition: if the Frobenius norm of  $(\Sigma - I)$  restricted on its  $(2k)^2$  largest entries (in magnitude) is small, then so is the maximum sparse eigenvalue – this would imply the empirical Chow vector be a good approximation to  $\chi_{f^*}$ ; see Theorem 17.

**2b) Large Frobenius norm validates a filter.** It remains to show that when the restricted Frobenius norm is large, say larger than some carefully chosen parameter  $\kappa$ , how to find a proper subset  $\bar{S}''$  that is clean enough, in the sense that its distributional properties act almost as it be an uncorrupted sample set. Our key idea is to construct a polynomial  $p_2$  such that (i) its empirical mean on the current sample set  $\bar{S}'$  equals the restricted Frobenius norm (which is large); and (ii) it has small value on clean samples. These two properties ensure that there must be a noticeable fraction of samples in  $\bar{S}'$  that caused a large function value of  $p_2$ , and they are very likely the corrupted ones; these will then be filtered to produce the new sample set  $\bar{S}''$  (Theorem 14). In addition, the polynomial  $p_2$  is constructed in such a way that it is *sparse* under the basis  $\{m_i(x) \cdot m_j(x)\}$ , for the sake of attribute efficiency.

We check the Frobenius norm condition every time a new sample set  $\bar{S}''$  is produced, and show that after a finite number of phases, we must be able to obtain a clean enough sample set  $\bar{S}''$  that allows us to output a good estimate of the Chow vector  $\chi_{f^*}$  (Theorem 18).

We note that the idea of using Frobenius norm as a surrogate of the maximum sparse eigenvalue value has been explored in Diakonikolas et al. (2019); Zeng & Shen (2022) for robust sparse mean estimation. In those works, the Frobenius-norm condition was combined with a localized eigenvalue condition to establish their main results, while we discover that the Frobenius norm itself suffices for our purpose. This appears an interesting and practical finding as it reduces the computational cost and simplifies algorithmic design.

## 1.3. Related works

The problem of learning from few samples dates back to the 1980s, when practitioners were confronting a pressing challenge: the number of samples available is orders of magnitude less than that of attributes, making classical algorithms fail to provide guarantees. The challenge persists

even in the big data era, since in many domains such as healthcare, there is a limited availability of samples (i.e. patients) (Candès & Wakin, 2008). This has motivated a flurry of research on attribute-efficient learning of sparse concepts. A partial list of interesting works includes (Littlestone, 1987; Blum, 1990; Chen et al., 1998; Tibshirani, 1996; Tropp, 2004; Candès & Tao, 2005; Foucart, 2011; Plan & Vershynin, 2013b; Shen & Li, 2018) that studied linear models in the absence of noise (or with benign noise). Later, Candès et al. (2013); Netrapalli et al. (2013); Candès et al. (2015) studied the problem of phase retrieval which can be seen as learning sparse quadratic polynomials. The setup was generalized in Chen & Meka (2020) which studied efficient learning of sparse low-degree polynomials.

In the presence of the nasty noise, the problem becomes subtle. Without distributional assumptions on  $D$ , it is known that even for the special case of learning halfspaces under the adversarial label noise, it is computationally hard when the noise rate is  $\epsilon$  (Guruswami & Raghavendra, 2006; Feldman et al., 2006; Daniely, 2016). Thus, distribution-independent algorithms are either unable to tolerate the nasty noise at a rate greater than  $\epsilon$  (Kearns & Li, 1988), or runs in super-polynomial time (Bshouty et al., 2002). This motivates the study of efficient algorithms under distributional assumptions (Kalai et al., 2005; Klivans et al., 2009; Awasthi et al., 2017; Shen & Zhang, 2021; Shen, 2023), which is the research line we follow. In unsupervised learning such as mean and covariance estimation, similar noise models are investigated broadly in recent years since the seminal works of Diakonikolas et al. (2016); Lai et al. (2016); see Diakonikolas & Kane (2019) for a comprehensive survey.

The interplay between sparsity and robustness is more involved than it appears to. Under the statistical-query framework, Diakonikolas et al. (2017) showed that any efficient and robust algorithms must draw  $\Omega(K^2)$  samples in the presence of the nasty noise, complementing sample complexity upper bounds obtained in recent years (Balakrishnan et al., 2017; Diakonikolas et al., 2019; Shen & Zhang, 2021; Diakonikolas et al., 2022). This is in stark contrast to learning with label noise, where  $O(K)$  sample complexity can be established (Zhang, 2018; Zhang et al., 2020; Shen, 2021a).

Lastly, we note that orthogonal to exploring sparsity for improved sample complexity, there are elegant works that explore sparsity for improved computational complexity for learning Boolean-valued functions (Hellerstein & Servedio, 2007; Andoni et al., 2014), or using low-degree PTFs as primitives to approximate other concepts such as halfspaces (Kalai et al., 2005) and decision lists (Servedio et al., 2012).

#### 1.4. Roadmap

We collect notations and definitions in Section 2. The main algorithms are described in Section 3 with a few lemmas to

illustrate the idea, and the primary performance guarantees are stated in Section 4. We conclude this work in Section 5. All omitted proofs can be found in the appendix.

## 2. Preliminaries

**Vectors and matrices.** For a vector  $v \in \mathbb{R}^n$ , we use  $v_i$  to denote its  $i$ -th element. For two vectors  $u$  and  $v$ , we write  $u \cdot v$  as the inner product in the Euclidean space. We denote by  $\|v\|_1$ ,  $\|v\|_2$ ,  $\|v\|_\infty$  the  $\ell_1$ -norm,  $\ell_2$ -norm, and  $\ell_\infty$ -norm of  $v$  respectively. The support set of a vector  $v$  is the index set of its non-zero elements, and  $\|v\|_0$  denotes the cardinality of the support set. We will use the hard thresholding operator  $H_k(v)$  to produce a  $k$ -sparse vector: the  $k$  largest (in magnitude) elements of  $v$  are retained and the rest are set to zero. Let  $\Lambda \subset [n]$  where  $[n] := \{1, \dots, n\}$ . The restriction of  $v$  on  $\Lambda$ ,  $v_\Lambda$ , is obtained by keeping the elements in  $\Lambda$  while setting the rest to zero.

Let  $A$  and  $B$  be two matrices in  $\mathbb{R}^{n_1 \times n_2}$ . We write  $\text{tr}(A)$  as the trace of  $A$  when it is square, and write  $\langle A, B \rangle := \text{tr}(A^\top B)$ . We denote by  $\|A\|_F$  the Frobenius norm, which equals  $\sqrt{\langle A, A \rangle}$ . We will also use  $\|A\|_0$  to count the number of non-zero entries in  $A$ . Let  $U \subset [n_1] \times [n_2]$ . The restriction of  $A$  on  $U$ ,  $A_U$ , is obtained by keeping the elements in  $U$  but setting the rest to zero.

**Probability,  $L^2$ -space.** Let  $D$  be a distribution on  $\mathbb{R}^n$  and  $p$  be a function with the same support of  $D$ . We denote by  $\mathbb{E}_{X \sim D}[p(X)]$  the expectation of  $p$  on  $D$ . Let  $S$  be a finite set of instances. We write  $\mathbb{E}_{X \sim S}[p(X)] := \frac{1}{|S|} \sum_{x \in S} p(x)$  as the empirical mean of  $p$  on  $S$ . To ease notation, we will often use  $\mathbb{E}[p(D)]$  in place of  $\mathbb{E}_{X \sim D}[p(X)]$ , and likewise for  $\mathbb{E}[p(S)]$ . Similarly, we will write  $\Pr(p(D) > t) := \Pr_{X \sim D}(p(X) > t)$ , and  $\Pr(p(S) > t) := \Pr_{X \sim S}(p(X) > t)$  where  $X \sim S$  signifies uniform distribution on  $S$ .

The  $L^2(\mathbb{R}^n, D)$  space is equipped with the inner product  $\langle p, q \rangle_D := \mathbb{E}_{x \sim D}[p(x) \cdot q(x)]$  for any functions  $p$  and  $q$  on  $\mathbb{R}^n$ . The induced  $L^2$ -norm of a function  $p$  is given by  $\|p\|_{L^2(D)} := \sqrt{\langle p, p \rangle_D} = \sqrt{\mathbb{E}[p^2(D)]}$ , which we will simply write as  $\|p\|_{L^2}$  when  $D$  is clear from the context.

**Polynomials.** Denote by  $\mathcal{P}_{n,d}$  the class of polynomials on  $\mathbb{R}^n$  with degree at most  $d$ . A  $\text{degree-}d$  *polynomial threshold function* (PTF) is of the form  $f(x) = \text{sign}(p(x))$  for some  $p \in \mathcal{P}_{n,d}$ . Denote by  $\text{He}_d(x) = \frac{1}{\sqrt{d!}} (-1)^d \cdot e^{-\frac{x^2}{2}} \frac{d^d}{dx^d} e^{-\frac{x^2}{2}}$  the normalized univariate degree- $d$  *Hermite polynomial* on  $\mathbb{R}$ . The normalized  $n$ -variate Hermite polynomial is given by  $\text{He}_\mathbf{a}(x) = \prod_{i=1}^n \text{He}_{a_i}(x_i)$  for some multi-index  $\mathbf{a} \in \mathbb{N}^n$ ; for brevity we refer to them as Hermite polynomials. It is known that  $\text{He}_{\leq d} := \{\text{He}_\mathbf{a} : \mathbf{a} \in \mathbb{N}^n, \|\mathbf{a}\|_1 \leq d\}$  form a complete orthonormal basis for polynomials of degree at most  $d$  in  $L^2(\mathbb{R}^n, D)$ ; see Prop. 11.33 of

O'Donnell (2014). It is easy to see that  $\text{He}_{\leq d}$  contains  $M := 1 + \sum_{t=1}^d \binom{t+n-1}{t}$  members; we collect them as a vector  $m(x) = (m_1(x), \dots, m_M(x))$ , with the first element  $m_1(x) \equiv 1$ . In our analysis, it suffices to keep in mind that  $M < (n+1)^d$ .

Given the vector  $m(x)$  and the distribution  $D$ , the *Chow vector* (Chow, 1961; Diakonikolas et al., 2018a) of a Boolean-valued function  $f : \mathbb{R}^n \rightarrow \{-1, 1\}$  is defined as follows:

$$\chi_f := \mathbb{E}_{x \sim D}[f(x) \cdot m(x)], \quad (2.1)$$

where we multiplied each element of  $m(x)$  by  $f(x)$ .

**Definition 9** (Sparse polynomials and PTFs). We say a polynomial  $p \in \mathcal{P}_{n,d}$  is  $K$ -sparse if there exists an index set  $\Lambda \subset [n]$  with  $|\Lambda| \leq K$ , such that  $p(x) = q(x_\Lambda)$  for some  $q \in \mathcal{P}_{n,d}$ . We say a PTF  $f(x) = \text{sign}(p(x))$  is  $K$ -sparse if  $p$  is  $K$ -sparse. The class of  $K$ -sparse PTFs on  $\mathbb{R}^n$  with degree at most  $d$  is denoted by  $\mathcal{H}_{d,K}$ .

One important observation is that our definition of sparse polynomials implies a sparsity pattern in the Chow vector; see Appendix B for the proof.

**Lemma 10.** Let  $f$  be a  $K$ -sparse degree- $d$  PTF. Then  $\chi_f$  is a  $k$ -sparse vector under the basis of Hermite polynomials, where  $k = d+1$  if  $K = 1$  and  $k \leq 2K^d$  otherwise.

As we discussed in Section I.2 there will be two concept classes involved in our algorithm and analysis. The first is the class of polynomials that have a sparse Chow vector under the basis of  $m(x)$ :

$$\mathcal{P}_{n,d,2k}^1 := \{p_1 : x \mapsto \langle v, m(x) \rangle, v \in \mathbb{R}^{(n+1)^d}, \|v\|_2 = 1, \|v\|_0 \leq 2k\}, \quad (2.2)$$

which will be useful in characterizing the approximation error to the Chow vector of interest. Another class consists of quadratic terms in  $m(x)$ ,

$$\mathcal{P}_{n,d,s}^2 := \{p_2 : x \mapsto \langle A_U, m(x)m(x)^\top - I \rangle, U^\top = U, \|U\|_0 \leq s, A \in \mathbb{S}^{(n+1)^d}, \|A_U\|_F = 1\} \quad (2.3)$$

where  $\mathbb{S}^{(n+1)^d} := \{A : A \in \mathbb{R}^{(n+1)^d \times (n+1)^d}, A^\top = A\}$ . Note that the polynomials in  $\mathcal{P}_{n,d,s}^2$  have degree at most  $2d$ , and can be represented as a linear combination of at most  $s$  elements of the form  $m_i(x)m_j(x)$ . They will be used to construct certain distributional statistics based on the empirical samples for filtering.

We will often use subscript to stress the membership of a polynomial in either class: we will write  $p_1 \in \mathcal{P}_{n,d,2k}^1$  and  $p_2 \in \mathcal{P}_{n,d,s}^2$ , rather than using the subscript for counting.

**Reserved symbols.** Throughout the paper,  $K$  always refers to the number of non-zero attributes that a sparse PTF depends on, and  $k$  is the sparsity of the Chow vector under

**Algorithm 1** Main Algorithm: Attribute-Efficient Robust Chow Vector Estimator

**Require:** A nasty adversary  $\text{EX}(\eta)$  with  $\eta \in [0, \frac{1}{2} - c]$  for some absolute constant  $c \in (0, \frac{1}{2}]$ , hypothesis class  $\mathcal{H}_{d,K}$  that contains  $f^*$ , target error rate  $\epsilon \in (0, 1)$ , failure probability  $\delta \in (0, 1)$ .

**Ensure:** A sparse vector  $u \in \mathbb{R}^{(n+1)^d}$ .

- 1:  $\bar{S}' \leftarrow \text{draw } C \cdot \frac{d^{5d} K^{4d}}{\epsilon^2} \log^{5d} \left( \frac{nd}{\epsilon\delta} \right) \text{ samples from } \text{EX}(\eta)$ .
- 2:  $k \leftarrow d+1$  if  $K = 1$  or  $k \leftarrow 2K^d$  if  $K > 1$ .
- 3:  $\kappa \leftarrow \frac{28}{c^2} \cdot [\rho_2 \cdot (c_0 \log \frac{1}{\eta} + c_0 d)^d \cdot \eta + \epsilon]$ .
- 4:  $l_{\max} \leftarrow \frac{4\eta k \gamma^2}{\epsilon} + 1$ .
- 5:  $\bar{S}'_l \leftarrow \bar{S}' \cap \{(x, y) : \|m(x)\|_\infty \leq \gamma\}$ .
- 6: **for** phase  $l = 1$  to  $l_{\max}$  **do**
- 7:    $\Sigma \leftarrow \mathbb{E}_{x \sim \bar{S}'_l} [m(x)m(x)^\top], \{(i_t, j_t)\}_{t=1}^{4k^2} \leftarrow \text{index set of the largest (in magnitude) } 2k \text{ diagonal entries and } 2k^2 - k \text{ entries above the main diagonal of } \Sigma - I$ .
- 8:    $U \leftarrow \{(i_t, j_t)\}_{t \geq 1} \cup \{(j_t, i_t)\}_{t \geq 1}$ .
- 9:   **if**  $\|\langle \Sigma - I, U \rangle\|_F \leq \kappa$  **then return**  $u \leftarrow H_k(\mathbb{E}_{(x,y) \sim \bar{S}'_l} [y \cdot m(x)])$ .
- 10: **end for**

**Algorithm 2** SPARSEFILTER( $S', U, \Sigma, k, \gamma, \rho_2$ )

- 1:  $A \leftarrow \frac{1}{\|\langle \Sigma - I, U \rangle\|_F} (\Sigma - I)_U$ .
- 2:  $p_2(x) \leftarrow \langle A, m(x)m(x)^\top - I \rangle$ .
- 3: Find  $t \in (0, 4k\gamma^2)$  such that  $\Pr(|p_2(S')| \geq t) \geq 6 \exp(- (t/\rho_2)^{1/d}/c_0) + \frac{3\epsilon}{k\gamma^2}$ .
- 4: **return**  $S'' \leftarrow \{x \in S' : |p_2(x)| \leq t\}$ .

the Hermite polynomials  $m(x)$  (see Lemma 10). We reserve  $\epsilon, \delta, \eta$  as in Definition 1. We wrote  $\bar{S}'$  as the corrupted sample set, and  $S'$  as the one without labels.

The capital and lowercase letters  $C$  and  $c$ , and their subscript variants, are always used to denote some absolute constants, though we do not track closely their values. We reserve

$$\gamma = (C_1 d \cdot \log \frac{nd}{\epsilon\delta})^{d/2}, \rho_2 = C_2 \cdot d^{\frac{3}{4}} \cdot (c_0 d)^d. \quad (2.4)$$

As will be clear in our analysis,  $\gamma$  upper bounds  $\max_{x \in S} \|m(x)\|_\infty$  for  $S$  drawn from  $D$ . Thus, we will only keep samples in  $\bar{S}'$  with  $x \in X_\gamma$  where

$$X_\gamma := \{x \in \mathbb{R}^n : \|m(x)\|_\infty \leq \gamma\}. \quad (2.5)$$

The quantity  $\rho_2$  upper bound  $\|p_2\|_{L^2}$ ; see Lemma 30.

### 3. Main Algorithms

Our main algorithm, Algorithm 1, aims to approximate the Chow vector of the underlying polynomial threshold

function  $f^* \in \mathcal{H}_{d,K}$  by drawing a small number of samples from the nasty adversary. Observe that the setting of  $k$  at Step 2 follows from Lemma 10, i.e.  $k$  is the sparsity of  $\chi_{f^*}$  under the basis of Hermite polynomials. With this in mind, we design three sparsity-induced components in the algorithm: pruning samples that must be outliers (Step 5), certifying that the sample set is clean enough and returning the empirical Chow vector (Step 8), or filtering samples with a carefully designed condition (Step 9). We elaborate on each component in the following.

### 3.1. Pruning

Since the outliers created by the nasty adversary may be arbitrary, it is useful to design some simple screening rule to remove samples that must have been corrupted. In this step, we leverage the distributional assumption that  $D$ , the distribution of instances, is the standard Gaussian  $\mathcal{N}(0, I_{n \times n})$ . Since the concept class consists of polynomials with degree at most  $d$ , it is known that any Hermite polynomial  $m_i(x)$  must concentrate around its mean with a known tail bound (Janson, 1997). As the mean of  $m_i(x)$  equals zero for all  $i \neq 1$  (recall that  $m_1(x) \equiv 1$ ), it is possible to specify a certain radius  $\gamma$  for pruning. Similar to Zeng & Shen (2022), we apply the  $\ell_\infty$ -norm metric for attribute efficiency, that is, we remove all samples  $(x, y)$  in  $\bar{S}'$  satisfying  $\|m(x)\|_\infty > \gamma$ . The following lemma shows that with high probability, no clean sample will be pruned under a proper choice of  $\gamma$ .

**Lemma 11.** *Let  $S$  be a set of samples drawn independently from  $D$ . With probability at least  $1 - \delta_\gamma$ , we have  $\max_{x \in S} \|m(x)\|_\infty \leq \gamma$  where  $\gamma := (c_0 \log \frac{|S|(n+1)^d}{\delta_\gamma})^{d/2}$ .*

Recall that the concrete value of  $\gamma$  is given in (2.4); it is obtained by setting  $|S|$  as the same size as in Step 1 of Algorithm 1 and setting  $\delta_\gamma = \frac{\epsilon^2 \delta}{64 \rho_2^2}$  (note that  $\delta_\gamma \leq O(\delta)$ ). The appearance of  $\delta$  in  $\delta_\gamma$  is not surprising. For the multiplicative factor  $\frac{\epsilon^2}{64 \rho_2^2}$ , technically speaking, it ensures that the total variation distance between the distribution  $D$  conditioned on the event  $x \in X_\gamma$  and  $D$  is  $O(\epsilon)$ , thus one can in principle consider uniform concentration on the former to ease analysis (since it is defined on a bounded domain); see Proposition 13.

### 3.2. Filtering

At Step 7 of Algorithm 1, we compute the empirical covariance matrix  $\Sigma$  and the index set  $U$  of the  $(2k)^2$  largest entries (in magnitude) of  $\Sigma - I$ . As we highlighted in Section 1.2, this is a computationally efficient way to obtain an upper bound on the maximum eigenvalue of  $\Sigma - I$  on all  $2k$ -sparse directions. The structural constraint on  $U$  comes from the observation that for  $2k$ -sparse  $v$ , we have  $v^\top (\Sigma - I)v = \langle \Sigma - I, vv^\top \rangle$  and  $vv^\top$  has  $2k$  non-zero diagonal entries and  $4k^2 - 2k$  off-diagonal symmetric entries.

If the restricted Frobenius norm,  $\|(\Sigma - I)_U\|_F$ , is greater than some threshold  $\kappa$ , Algorithm 1 will invoke a filtering subroutine, Algorithm 2 to remove samples that were potentially corrupted. The high-level idea of Algorithm 2 follows from prior works on robust mean estimation (Diakonikolas et al., 2016; 2019; Zeng & Shen, 2022): under the condition that a certain measure of the empirical covariance matrix is large, there must be some samples that behave in quite a different way from those drawing from  $D$ . Our technical novelty is a new algorithm and analysis showing that the Frobenius norm itself suffices to validate a certain type of test that can identify those potentially corrupted samples – this is a new feature as existing robust sparse mean estimation algorithms (Diakonikolas et al., 2019; Zeng & Shen, 2022) rely on a combination of the Frobenius norm and a localized eigenvalue condition. An immediate implication of our finding is that one can expect lower computational cost of our algorithm due to the lack of eigenvalue computation.

Now we discuss how to design a test to filter potentially corrupted samples. The idea is to create a sample-dependent polynomial  $p_2$  with the following two properties: 1) its empirical mean on  $S'$  equals  $\|(\Sigma - I)_U\|_F$ ; and 2)  $p_2$  is small (in expectation) on uncorrupted samples. In this way, since we have the condition that  $\|(\Sigma - I)_U\|_F$  is large, there must be a noticeable fraction of samples in  $S'$  that correspond to large function values of  $p_2$ . This combined with the second property suffice to identify abnormal samples.

Indeed, since  $\Sigma = \mathbb{E}_{x \sim S'}[m(x)m(x)^\top]$ , we can show that

$$\begin{aligned} & \|(\Sigma - I)_U\|_F \\ &= \frac{1}{\|(\Sigma - I)_U\|_F} \langle (\Sigma - I)_U, \mathbb{E}_{x \sim S'}[m(x)m(x)^\top] - I \rangle \\ &= \mathbb{E}_{x \sim S'} \left[ \frac{1}{\|(\Sigma - I)_U\|_F} \langle (\Sigma - I)_U, m(x)m(x)^\top - I \rangle \right]. \end{aligned}$$

This gives the design of  $p_2$  in Algorithm 2 which has the desired feature: its expectation on  $D$  equals zero since  $m(x)$  is an orthonormal basis in  $L^2(\mathbb{R}^n, D)$ . Yet, we remark that the degree of  $p_2$  is as large as  $2d$ , which leads to a heavy-tailed distribution even for uncorrupted data; and thus the nasty adversary may inject comparably heavy-tailed data. In Lemma 30, we show that the  $L^2$ -norm of  $p_2$  on  $D$  is upper bounded by  $\rho_2 = O(d^d)$ ; thus the threshold  $\kappa$  is proportional to  $\rho_2$ . The additional multiplicative factor in  $\kappa$ ,  $(c_0 \log \frac{1}{\eta} + c_0 d)^d \cdot \eta$ , is the maximum amount that those  $\eta$ -fraction of heavy-tailed outliers can deteriorate the restricted Frobenius norm without appearing quite different from uncorrupted samples. In other words, with this scaling of  $\kappa$ , if the outliers were to deviate our estimate significantly, they would trigger the filtering condition.

Now we give intuition on Step 3 of Algorithm 2. We can use standard results on Gaussian tail bound of polynomials

(Janson, 1997) to show that

$$\Pr(|p_2(D)| \geq t) \leq \exp(-(t/\rho_2)^{1/d}/c_0), \forall t > 0.$$

By uniform convergence (Vapnik & Chervonenkis, 1971), the above implies a low frequency of the event  $|p_2(x)| \geq t$  on a set of uncorrupted samples (provided that the sample size is large enough; see Part 5 of Definition 12). On the other hand, the empirical average of  $p_2$  on the input instance set  $S'$  (which equals  $\|(\Sigma - I)_{U'}\|_F$ ) is large. Thus, there must be some threshold  $t$  such that  $|p_2(x)| \geq t$  occurs with a nontrivial frequency, and this is an indicator of being outliers. In Step 3 of Algorithm 2, the nontrivial frequency is set as a constant factor of the one of uncorrupted samples – it is known that this suffices to produce a cleaner instance set; see e.g. Diakonikolas et al. (2016). To further guarantee a bounded running time, we show that it suffices to find a  $t$  in  $(0, 4k\gamma^2)$ , thanks to the pruning step (see Lemma 30).

It is worth mentioning that our primary treatment on attribute efficiency lies in applying uniform convergence to derive the low frequency event. In fact, since the size of  $U$  is at most  $4k^2$ , it is possible to show that the VC-dimension of the class  $\mathcal{P}_{n,d,s}^2$  that  $p_2$  resides is  $O(s \log n^d)$ , with  $s = 4k^2$ .

### 3.3. Termination

Lastly, we describe the case that Algorithm 1 terminates and output  $u$  at Step 8. Due to the selection of  $U$ , it is possible to show that  $\|(\Sigma - I)_{U'}\|_F \leq \kappa$  implies  $v^\top \Sigma v \leq \kappa + 1$  for all  $2k$ -sparse unit vector  $v$ , i.e. the maximum eigenvalue of  $\Sigma$  on all  $2k$ -sparse directions is as small as  $\kappa + 1$ . This in turn implies that the variance caused by corrupted samples is well-controlled. Therefore, we output the empirical Chow vector. We note that Algorithm 1 outputs  $u$  which is the empirical one followed by a hard thresholding operation. This ensures that  $u$  is  $k$ -sparse, the same sparsity level as  $\chi_{f^*}$ . More importantly, since we are only guaranteed with a small maximum eigenvalue on  $2k$ -sparse directions, it is likely that on the full direction, the maximum eigenvalue could be very large, which would fail to certify a good approximation to  $\chi_{f^*}$ . In other words, had we not applied the hard thresholding operation, the empirical estimate  $\mathbb{E}_{(x,y) \sim S'_l} [y \cdot m(x)]$  could be far away from the target Chow vector.

The maximum number of iterations,  $l_{\max}$ , comes from our analysis on the progress of the filtering step: we will show in Section 4 that each time Algorithm 2 is invoked, a noticeable fraction of outliers will be removed while most clean samples are retained, thus after at most  $l_{\max}$  iterations, the restricted Frobenius norm must be less than  $\kappa$ .

## 4. Performance Guarantees

Our analysis of filtering relies on the existence of a good set  $S \subset \mathbb{R}^n$  and shows that Algorithm 2 strictly reduces

the distance between the corrupted set and  $S$  every time it is invoked by Algorithm 1, until the termination condition is met (Theorem 14). We then show that the output of Algorithm 1 must be close to the Chow vector of the underlying PTF (Theorem 17), and this occurs within  $l_{\max}$  phases (Theorem 18). Then, a black-box application of the algorithmic result from Trevisan et al. (2009); De et al. (2014); Diakonikolas et al. (2018a) leads to PAC guarantees of a PTF that is reconstructed from our estimated Chow vector (Theorem 19).

We will phrase our results in terms of some deterministic conditions on  $S$ . Let  $S|_{X_\gamma} := S \cap X_\gamma$  and  $D|_{X_\gamma}$  be the distribution  $D$  conditioned on the event  $x \in X_\gamma$ .

**Definition 12** (Good set). Given  $\epsilon \in (0, 1)$ ,  $\delta \in (0, 1)$ , and concept class  $\mathcal{H}_{d,K}$ , we say an instance set  $S \subset \mathbb{R}^n$  is a good set if all the following properties hold simultaneously and uniformly over all  $p_1 \in \mathcal{P}_{n,d,2k}^1$  ( $k$  is given in Lemma 10), all  $p_2 \in \mathcal{P}_{n,d,s}^2$  with  $s = 4k^2$ , and all  $t > 0$ :

1.  $S|_{X_\gamma} = S$ ;
2.  $|\Pr(p_1(S) > t) - \Pr(p_1(D) > t)| \leq \alpha_1$ ;
3.  $|\Pr(p_1(S|_{X_\gamma}) > t) - \Pr(p_1(D|_{X_\gamma}) > t)| \leq \alpha_1$ ;
4.  $|\mathbb{E}_{x \sim S} [f(x) \cdot p_1(x)] - \mathbb{E}_{x \sim D} [f(x) \cdot p_1(x)]| \leq \alpha'_1$ ;
5.  $|\Pr(p_2(S) > t) - \Pr(p_2(D) > t)| \leq \alpha_2$ ;
6.  $|\Pr(p_2(S|_{X_\gamma}) > t) - \Pr(p_2(D|_{X_\gamma}) > t)| \leq \alpha_2$ ;
7.  $|\mathbb{E}[p_2(S)] - \mathbb{E}[p_2(D)]| \leq \alpha'_2$ ,

where  $\alpha_1 = \frac{\epsilon}{k\gamma^2}$ ,  $\alpha'_1 = \epsilon/6$ ,  $\alpha_2 = \frac{\epsilon}{4k\gamma^2}$ ,  $\alpha'_2 = \epsilon$ .

We show that for a set of instances independently drawn from  $D$ , it is indeed a good set. Note that this gives the sample size at Step 1 of Algorithm 1.

**Proposition 13.** Let  $S$  be a set of  $C \cdot \frac{d^{5d} K^{4d}}{\epsilon^2} \log^{5d} \left( \frac{nd}{\epsilon\delta} \right)$  instances drawn independently from  $D$ . Then with probability  $1 - \delta$ ,  $S$  is a good set.

### 4.1. Analysis of SPARSEFILTER

Recall in Definition 1 that the nasty adversary first draws  $S$  according to  $D$  and annotates it with  $f^*$  to obtain  $\bar{S} \subset \mathbb{R}^n \times \{-1, 1\}$ . Then it replaces an  $\eta$  fraction with malicious samples to generate the sample set  $\bar{S}'$  that is returned to the learner. Denote by  $\Delta(S, S')$  the symmetric difference between  $S$  and  $S'$  normalized by  $|S|$ , i.e.

$$\Delta(S, S') := \frac{|S \setminus S'| + |S' \setminus S|}{|S|}. \quad (4.1)$$

By definition, it follows that  $\Delta(S, S') \leq 2\eta$ . The following theorem is the primary characterization of the performance of our filtering approach (Algorithm 2).

**Theorem 14.** *Consider Algorithm 2. Assume that  $\|(\Sigma - I)_U\|_F > \kappa$  and there exists a good set  $S$  such that  $\Delta(S, S') \leq 2\eta$ . Then there exists a threshold  $t$  that satisfies Step 3. In addition, the output  $S''$  satisfies  $\Delta(S, S'') \leq \Delta(S, S') - \frac{\epsilon}{4k\gamma^2}$ .*

We show this theorem by contradiction: had we been unable to find such  $t$ , the tail bound at Step 3 would have implied small expectation of  $p_2$  on  $S'$ . As discussed in Section 3.2 the polynomial  $p_2$  is chosen such that  $\|(\Sigma - I)_U\|_F = \mathbb{E}[p_2(S')]$ ; this in turn suggests that we would contradict the condition  $\|(\Sigma - I)_U\|_F > \kappa$  when  $\kappa$  is properly chosen.

Formally, let  $\beta'(\tau, d, \rho) := 2 \cdot \rho \cdot (c_0 \log \frac{1}{\tau} + c_0 \cdot d/2)^{d/2} \cdot \tau$  and  $\gamma_2 := 4k^2\gamma^2$ . We have:

**Lemma 15.** *Consider Algorithm 2. Assume that  $\|(\Sigma - I)_U\|_F > \kappa$  and there exists a good set  $S$  such that  $\Delta(S, S') \leq 2\eta$ . Let  $E := S' \setminus S$ . If there does not exist a threshold  $t > 0$  that satisfies Step 3, then*

$$\frac{|E|}{|S'|} \sup_{p_2 \in \mathcal{P}_{n,d,s}^2} \mathbb{E}[|p_2(E)|] \leq 7(1 + \frac{1}{c}) \cdot [\beta'(\eta, 2d, \rho_2) + \alpha_2 \gamma_2].$$

**Lemma 16.** *Consider Algorithm 2. Assume that there exists a good set  $S$  with  $\Delta(S, S') \leq 2\eta$ . Let  $L := S \setminus S'$ . We have*

$$\frac{|L|}{|S|} \sup_{p_2 \in \mathcal{P}_{n,d,s}^2} \mathbb{E}[|p_2(L)|] \leq 2(1 + \frac{1}{c}) [\beta'(\eta, 2d, \rho_2) + \alpha_2 \gamma_2].$$

Now observe that  $|S'| \cdot \|(\Sigma - I)_U\|_F = |S'| \cdot \mathbb{E}[p_2(S')] = |S| \cdot \mathbb{E}[p_2(S)] + |E| \cdot \mathbb{E}[p_2(E)] - |L| \cdot \mathbb{E}[p_2(L)]$ . For the right-hand side, we can roughly think of  $\mathbb{E}[p_2(S)] \approx \mathbb{E}[p_2(D)]$  which can be bounded as  $D$  is Gaussian. This combined with Lemma 15 and Lemma 16 can establish the existence of  $t$ . We then use a general result that is implicit in prior filter-based algorithms (Diakonikolas et al., 2016): given the existence of  $t$ , there must be a nontrivial fraction of the instances in  $S'$  that can be filtered; see also Lemma 38 where we provide a generic proof. This establishes Theorem 14, see Appendix D for the full proof.

## 4.2. Analysis of termination

Let  $\beta_\tau = 2(c_0 \log \frac{1}{\tau} + c_0 d)^{d/2} \cdot \tau$  for some parameter  $\tau \in (0, 1)$ . The following theorem shows that whenever the termination condition is met, i.e.  $\|(\Sigma - I)_U\|_F \leq \kappa$ , the output must be close to the target Chow vector.

**Theorem 17.** *Consider Algorithm 1. If at some phase  $l$  we have  $\|(\Sigma - I)_U\|_F \leq \kappa$  and  $\Delta(S, S'_l) \leq 2\eta$  for some good set  $S$ , then the following holds for the output  $u$ :*

$$\|u - \chi_{f^*}\|_2 \leq \frac{192}{c^2} \sqrt{\eta(\beta_\eta + \beta_\epsilon)} + \frac{\epsilon}{2}.$$

We note that the upper bound seems not depending on  $\kappa$  – this is because  $\kappa \leq O(\beta_\epsilon)$ . To show the theorem, we will first prove that the deviation of the expectation of  $y \cdot p_1(x)$  between  $\bar{S}'$  and  $\bar{S}$  is small, and then apply Part 4 of Definition 12 to establish the closeness to the expectation on  $D$ . To obtain the first deviation bound, we observe that it is almost governed by the expectation on  $S \setminus S'$  and on  $S' \setminus S$ . The former is easy to control since it is a subset of the good set  $S$ . We show that the latter is also bounded since the termination condition implies a small variance on all sparse directions of the covariance matrix  $\Sigma$  that is computed on  $S'$ ; this suggests that the contribution from the corrupted instances cannot be large. See Appendix E for the proof.

## 4.3. Main results

**Theorem 18** (Chow vector estimation). *The following holds for Algorithm 1. Given any target error rate  $\epsilon \in (0, 1)$  and failure probability  $\delta \in (0, 1)$ , Algorithm 1 runs in at most  $l_{\max} = \frac{4\eta k}{\epsilon} \cdot (C_1 d \cdot \log \frac{nd}{\epsilon\delta})^d + 1$  phases, and outputs a  $k$ -sparse vector  $u$  such that with probability at least  $1 - \delta$ ,*

$$\|u - \chi_{f^*}\|_2 \leq \frac{192}{c^2} \sqrt{\eta(\beta_\eta + \beta_\epsilon)} + \frac{\epsilon}{2}.$$

*In addition, Algorithm 1 runs in  $O(\text{poly}((nd)^d, 1/\epsilon))$  time.*

*Proof sketch.* In view of Proposition 13 and Step 1 of Algorithm 1, there is a good set  $S$  such that  $\Delta(S, S') \leq 2\eta$ . We will inductively show the invariant  $\Delta(S, S'_{l+1}) \leq \Delta(S, S'_l) - \frac{\epsilon}{4k\gamma^2}$  before Algorithm 1 terminates. In fact, by Part 1 of Definition 12, it follows that no instances in  $S$  will be pruned at Step 1. Thus,  $\Delta(S, S'_1) \leq \Delta(S, S') \leq 2\eta$ . If  $\|(\Sigma - I)_U\|_F > \kappa$ , then Theorem 14 implies that we can obtain  $S'_2$  such that  $\Delta(S, S'_2) \leq \Delta(S, S'_1) - \frac{\epsilon}{4k\gamma^2}$ . By induction, we can show that such progress holds for any phase  $l$  before the termination condition is met. Since the symmetric difference is non-negative, the algorithm must terminate within the claimed  $l_{\max}$  phases, upon when the output is guaranteed to be close to  $\chi_{f^*}$  in view of Theorem 17. See Appendix E for the full proof.  $\square$

Lastly, the algorithmic results from Trevisan et al. (2009); De et al. (2014); Diakonikolas et al. (2018a) state that as long as  $u$  is  $\epsilon$ -close to  $\chi_{f^*}$  under the  $\ell_2$ -norm, it is possible to construct a PTF  $\hat{f}$  in time  $(\frac{n}{\epsilon})^{O(d)}$  that has misclassification error of  $O(d \cdot \epsilon^{1/(d+1)})$ . This gives our main result on PAC guarantees (see Appendix F for the proof).

**Theorem 19** (PAC guarantees). *There exists an algorithm  $\mathcal{A}$  such that the following holds. Given any  $\epsilon_0 \in (0, 1)$ , failure probability  $\delta \in (0, 1)$ , and the concept class  $\mathcal{H}_{d,K}$ , it draws  $C \cdot \frac{d^{5d} K^{4d}}{\epsilon_0^2} \cdot \log^{5d} (\frac{nd}{\epsilon_0 \delta})$  samples from  $\text{EX}(\eta)$  and*

outputs a PTF  $\hat{f}$  such that with probability at least  $1 - \delta$ ,  $\Pr_{x \sim D}(\hat{f}(x) \neq f^*(x)) \leq c_1 \cdot d \cdot \left( \sqrt{\eta(\beta_\eta + \beta_{\epsilon_0})} + \epsilon_0 \right)^{\frac{1}{d+1}}$ .

In particular, for any target error rate  $\epsilon \in (0, 1)$ , by setting  $\epsilon_0 = \frac{\epsilon^{d+1}}{c_2 \cdot d^2}$ , we have  $\Pr_{x \sim D}(\hat{f}(x) \neq f^*(x)) \leq \epsilon$  provided  $\eta \leq \frac{1}{2}\epsilon_0$ . Moreover, the algorithm runs in  $(nd/\epsilon)^{O(d)}$  time.

## 5. Conclusion

We studied the important problem of attribute-efficient PAC learning of low-degree PTFs. We showed that for the class of sparse PTFs where the concept depends only on a subset of its input attributes, it is possible to design an efficient algorithm that PAC learns the class with sample complexity poly-logarithmic in the dimension, even in the presence of the nasty noise. In addition, the noise tolerance of our algorithm is dimension-independent, and matches the best known result established for learning of non-sparse PTFs.

## Acknowledgements

We thank the anonymous reviewers for valuable comments. We thank Aughdon Breslin and Matthew Thomas for helping calculate the sparsity of the Chow vector under monomials. This work is supported by NSF-AF-2239376 and the startup funding from Stevens Institute of Technology.

## References

Andoni, A., Panigrahy, R., Valiant, G., and Zhang, L. Learning sparse polynomial functions. In *Proceedings of the 25th Annual ACM-SIAM Symposium on Discrete Algorithms*, pp. 500–510, 2014.

Anthony, M. and Bartlett, P. L. *Neural Network Learning: Theoretical Foundations*. Cambridge University Press, 1999.

Awasthi, P., Balcan, M., Haghtalab, N., and Zhang, H. Learning and 1-bit compressed sensing under asymmetric noise. In *Proceedings of the 29th Annual Conference on Learning Theory*, pp. 152–192, 2016.

Awasthi, P., Balcan, M., and Long, P. M. The power of localization for efficiently learning linear separators with noise. *Journal of the ACM*, 63(6):50:1–50:27, 2017.

Balakrishnan, S., Du, S. S., Li, J., and Singh, A. Computationally efficient robust sparse estimation in high dimensions. In *Proceedings of the 30th Annual Conference on Learning Theory*, pp. 169–212, 2017.

Blum, A. Learning boolean functions in an infinite attribute space. In *Proceedings of the 22nd Annual ACM Symposium on Theory of Computing*, pp. 64–72, 1990.

Blumer, A., Ehrenfeucht, A., Haussler, D., and Warmuth, M. K. Learnability and the Vapnik-Chervonenkis dimension. *Journal of the ACM*, 36(4):929–965, 1989.

Bshouty, N. H., Eiron, N., and Kushilevitz, E. PAC learning with nasty noise. *Theoretical Computer Science*, 288(2): 255–275, 2002.

Candès, E. J. and Tao, T. Decoding by linear programming. *IEEE Transactions on Information Theory*, 51(12):4203–4215, 2005.

Candès, E. J. and Wakin, M. B. An introduction to compressive sampling. *IEEE Signal Processing Magazine*, 25(2): 21–30, 2008.

Candès, E. J., Strohmer, T., and Voroninski, V. Phaselift: Exact and stable signal recovery from magnitude measurements via convex programming. *Communications on Pure and Applied Mathematics*, 66(8):1241–1274, 2013.

Candès, E. J., Li, X., and Soltanolkotabi, M. Phase retrieval via Wirtinger flow: Theory and algorithms. *IEEE Transactions on Information Theory*, 61(4):1985–2007, 2015.

Chen, S. and Meka, R. Learning polynomials in few relevant dimensions. In *Proceedings of the 34th Annual Conference on Learning Theory*, pp. 1161–1227, 2020.

Chen, S. S., Donoho, D. L., and Saunders, M. A. Atomic decomposition by basis pursuit. *SIAM Journal on Scientific Computing*, 20(1):33–61, 1998.

Chow, C.-K. On the characterization of threshold functions. In *Proceedings of the 2nd Annual Symposium on Switching Circuit Theory and Logical Design (FOCS)*, pp. 34–38, 1961.

Daniely, A. Complexity theoretic limitations on learning halfspaces. In *Proceedings of the 48th Annual ACM Symposium on Theory of Computing*, pp. 105–117, 2016.

De, A., Diakonikolas, I., Feldman, V., and Servedio, R. A. Nearly optimal solutions for the Chow parameters problem and low-weight approximation of halfspaces. *Journal of the ACM*, 61(2):11:1–11:36, 2014.

Diakonikolas, I. and Kane, D. M. Recent advances in algorithmic high-dimensional robust statistics. *CoRR*, abs/1911.05911, 2019.

Diakonikolas, I., Kamath, G., Kane, D. M., Li, J., Moitra, A., and Stewart, A. Robust estimators in high dimensions without the computational intractability. In *Proceedings of the 57th Annual IEEE Symposium on Foundations of Computer Science*, pp. 655–664, 2016.

Diakonikolas, I., Kane, D. M., and Stewart, A. Statistical query lower bounds for robust estimation of high-dimensional gaussians and gaussian mixtures. In *Proceedings of the 58th IEEE Annual Symposium on Foundations of Computer Science*, pp. 73–84, 2017.

Diakonikolas, I., Kane, D. M., and Stewart, A. Learning geometric concepts with nasty noise. In *Proceedings of the 50th Annual ACM Symposium on Theory of Computing*, pp. 1061–1073, 2018a.

Diakonikolas, I., Kane, D. M., and Stewart, A. List-decodable robust mean estimation and learning mixtures of spherical gaussians. In *Proceedings of the 50th Annual ACM SIGACT Symposium on Theory of Computing*, pp. 1047–1060, 2018b.

Diakonikolas, I., Kane, D., Karmalkar, S., Price, E., and Stewart, A. Outlier-robust high-dimensional sparse estimation via iterative filtering. In *Proceedings of the 33rd Annual Conference on Neural Information Processing Systems*, pp. 10688–10699, 2019.

Diakonikolas, I., Kane, D. M., Karmalkar, S., Pensia, A., and Pittas, T. Robust sparse mean estimation via sum of squares. In *Proceedings of the The 35th Annual Conference on Learning Theory*, pp. 4703–4763, 2022.

Feldman, V., Gopalan, P., Khot, S., and Ponnuswami, A. K. New results for learning noisy parities and halfspaces. In *Proceedings of the 47th Annual IEEE Symposium on Foundations of Computer Science*, pp. 563–574, 2006.

Foucart, S. Hard thresholding pursuit: An algorithm for compressive sensing. *SIAM Journal on Numerical Analysis*, 49(6):2543–2563, 2011.

Gentile, C. The robustness of the  $p$ -norm algorithms. *Machine Learning*, 53(3):265–299, 2003.

Guruswami, V. and Raghavendra, P. Hardness of learning halfspaces with noise. In *Proceedings of the 47th Annual IEEE Symposium on Foundations of Computer Science*, pp. 543–552, 2006.

Hellerstein, L. and Servedio, R. A. On PAC learning algorithms for rich Boolean function classes. *Theoretical Computer Science*, 384(1):66–76, 2007.

Janson, S. *Gaussian Hilbert Spaces*. Cambridge Tracts in Mathematics. Cambridge University Press, 1997.

Kalai, A. T., Klivans, A. R., Mansour, Y., and Servedio, R. A. Agnostically learning halfspaces. In *Proceedings of the 46th Annual IEEE Symposium on Foundations of Computer Science*, pp. 11–20, 2005.

Kearns, M. J. and Li, M. Learning in the presence of malicious errors. In *Proceedings of the 20th Annual ACM Symposium on Theory of Computing*, pp. 267–280, 1988.

Klivans, A. R., Long, P. M., and Servedio, R. A. Learning halfspaces with malicious noise. *Journal of Machine Learning Research*, 10:2715–2740, 2009.

Lai, K. A., Rao, A. B., and Vempala, S. S. Agnostic estimation of mean and covariance. In *Proceedings of the 57th Annual IEEE Symposium on Foundations of Computer Science*, pp. 665–674, 2016.

Littlestone, N. Learning quickly when irrelevant attributes abound: A new linear-threshold algorithm. In *Proceedings of the 28th Annual IEEE Symposium on Foundations of Computer Science*, pp. 68–77, 1987.

Long, P. M. and Servedio, R. A. Learning large-margin halfspaces with more malicious noise. In *Proceedings of the 25th Annual Conference on Neural Information Processing Systems*, pp. 91–99, 2011.

Maass, W. and Turán, G. How fast can a threshold gate learn? In *Proceedings of a workshop on computational learning theory and natural learning systems (vol. 1): constraints and prospects*, pp. 381–414, 1994.

Mansour, Y. Learning boolean functions via the Fourier transform. In *Theoretical advances in neural computation and learning*, pp. 391–424. Springer, 1994.

Netrapalli, P., Jain, P., and Sanghavi, S. Phase retrieval using alternating minimization. In *Proceedings of the 27th Annual Conference on Neural Information Processing Systems*, pp. 2796–2804, 2013.

O’Donnell, R. *Analysis of Boolean Functions*. Cambridge University Press, 2014.

Plan, Y. and Vershynin, R. One-bit compressed sensing by linear programming. *Communications on Pure and Applied Mathematics*, 66(8):1275–1297, 2013a.

Plan, Y. and Vershynin, R. Robust 1-bit compressed sensing and sparse logistic regression: A convex programming approach. *IEEE Transactions on Information Theory*, 59(1):482–494, 2013b.

Servedio, R. A., Tan, L., and Thaler, J. Attribute-efficient learning and weight-degree tradeoffs for polynomial threshold functions. In *Proceedings of the 25th Annual Conference on Learning Theory*, pp. 1–19, 2012.

Shen, J. On the power of localized Perceptron for label-optimal learning of halfspaces with adversarial noise. In *Proceedings of the 38th International Conference on Machine Learning*, pp. 9503–9514, 2021a.

Shen, J. Sample-optimal PAC learning of halfspaces with malicious noise. In *Proceedings of the 38th International Conference on Machine Learning*, pp. 9515–9524, 2021b.

Shen, J. PAC learning of halfspaces with malicious noise in nearly linear time. In *Proceedings of the 26th International Conference on Artificial Intelligence and Statistics*, pp. 30–46, 2023.

Shen, J. and Li, P. A tight bound of hard thresholding. *Journal of Machine Learning Research*, 18(208):1–42, 2018.

Shen, J. and Zhang, C. Attribute-efficient learning of halfspaces with malicious noise: Near-optimal label complexity and noise tolerance. In *Proceedings of the 32nd International Conference on Algorithmic Learning Theory*, pp. 1072–1113, 2021.

Tibshirani, R. Regression shrinkage and selection via the Lasso. *Journal of the Royal Statistical Society: Series B (Methodological)*, 58(1):267–288, 1996.

Trevisan, L., Tulsiani, M., and Vadhan, S. P. Regularity, boosting, and efficiently simulating every high-entropy distribution. In *Proceedings of the 24th Annual IEEE Conference on Computational Complexity*, pp. 126–136, 2009.

Tropp, J. A. Greed is good: algorithmic results for sparse approximation. *IEEE Transactions on Information Theory*, 50(10):2231–2242, 2004.

Valiant, L. G. A theory of the learnable. *Communications of the ACM*, 27(11):1134–1142, 1984.

Vapnik, V. N. and Chervonenkis, A. Y. On the uniform convergence of relative frequencies of events to their probabilities. *Theory of Probability and its Applications*, 16(2):264, 1971.

Zeng, S. and Shen, J. List-decodable sparse mean estimation. In *Proceedings of the 36th Annual Conference on Neural Information Processing Systems*, pp. 24031–24045, 2022.

Zhang, C. Efficient active learning of sparse halfspaces. In *Proceedings of the 31st Annual Conference On Learning Theory*, pp. 1856–1880, 2018.

Zhang, C., Shen, J., and Awasthi, P. Efficient active learning of sparse halfspaces with arbitrary bounded noise. In *Proceedings of the 34th Annual Conference on Neural Information Processing Systems*, pp. 7184–7197, 2020.

We summarize a few useful results and list reserved hyper-parameters in Appendix A, they will be frequently used in our analysis. We provide proofs for results in Section 1 and Section 2 in Appendix B. Appendix C collects statistical results on the concept classes of interest, which are used in Appendix D and Appendix E to establish guarantees on Algorithm 2 and Algorithm 1, respectively. We assemble all pieces and prove the main result, Theorem 19 in Appendix F.

## A. Summary of Useful Facts and Reserved Hyper-Parameters

We will often need the condition that  $\Delta(S, S') \leq 2\eta$ , which implies

$$(1 - 2\eta) |S| \leq |S'| \leq |S|. \quad (\text{A.1})$$

In particular, when  $\eta \in [0, \frac{1}{2} - c]$  for some absolute constant  $c \in (0, \frac{1}{2}]$ , we have

$$|S'| \leq |S| \leq \frac{1}{1 - 2\eta} |S'| \leq \left(1 + \frac{1}{c} \cdot \eta\right) |S'|. \quad (\text{A.2})$$

The above two inequalities also imply

$$|S' \setminus S| \leq 2\eta |S| \leq \frac{2\eta}{1 - 2\eta} |S'| \leq \frac{\eta}{c} |S'| \quad \text{and} \quad |S \setminus S'| \leq 2\eta |S| \leq \frac{2\eta}{1 - 2\eta} |S'| \leq \frac{\eta}{c} |S'|. \quad (\text{A.3})$$

It is known that for any vector  $u$ ,

$$\|u\|_1 \leq \sqrt{\|u\|_0 \cdot \|u\|_2}. \quad (\text{A.4})$$

The above will often be applied together with Holder's inequality:

$$|u \cdot v| \leq \|u\|_1 \cdot \|v\|_\infty \leq \sqrt{\|u\|_0 \cdot \|u\|_2} \cdot \|v\|_\infty. \quad (\text{A.5})$$

**Fact 20.** Let  $Z$  be a positive random variable. Then  $\mathbb{E}[Z] = \int_0^\infty \Pr(Z > t) dt$ .

**Fact 21** (Tail bound of Gaussian polynomials (Janson, 1997)). Let  $D$  be the standard Gaussian distribution  $\mathcal{N}(0, I_{n \times n})$ . There exists an absolute constant  $c_0 > 1$  such that the following tail bound holds for all degree- $d$  polynomials  $p : \mathbb{R}^n \rightarrow \mathbb{R}$ :

$$\Pr_{x \sim D} (|p(x) - \mathbb{E}[p(D)]| \geq t \sqrt{\text{Var}[p(D)]}) \leq \exp(-t^{2/d}/c_0), \quad \forall t > 0.$$

In particular, if  $p$  is such that  $\mathbb{E}[p(D)] = 0$ , we have

$$\Pr_{x \sim D} (|p(x)| \geq t \|p\|_{L^2}) \leq \exp(-t^{2/d}/c_0).$$

### A.1. Reserved Hyper-Parameters

Recall that  $\epsilon \in (0, 1)$  is the noise rate,  $\delta \in (0, 1)$  is the failure probability,  $d$  is the degree of the PTFs. Denote by  $X_\gamma = \{x \in \mathbb{R}^n : \|m(x)\|_\infty \leq \gamma\}$  the instances of interest. Given an instance set  $S \subset \mathbb{R}^n$ , let  $S|_{X_\gamma} = S \cap X_\gamma$ . For a distribution  $D$  supported on  $\mathbb{R}^n$ , let  $D|_{X_\gamma}$  be the distribution  $D$  conditioned on the event that  $x \in X_\gamma$ .

- $\beta(\tau, d, \rho) = \rho^2 \cdot (c_0 \log \frac{1}{\tau} + c_0 d)^d \cdot \tau$ , which upper bounds  $\int_0^\infty t \cdot \min\{\tau, Q_{\rho, d}(t)\} dt$  for  $Q_{d, \rho}(t) = \exp(-(t/\rho)^{2/d}/c_0)$ ; see Lemma 24;
- $\beta'(\tau, d, \rho) = 2 \cdot \rho \cdot (c_0 \log \frac{1}{\tau} + c_0 \cdot d/2)^{d/2} \cdot \tau$ , which upper bounds  $\int_0^\infty \min\{\tau, Q_{\rho, d}(t)\} dt$ ; see Lemma 27;
- $\gamma = (c_0 \log \frac{|S| \cdot (n+1)^d}{\delta_\gamma})^{d/2} = (C_1 d \cdot \log \frac{nd}{\epsilon \delta})^{d/2}$ , which upper bounds  $\max_{x \in S} \|m(x)\|_\infty$  with probability  $1 - \delta_\gamma$  for  $S$  drawn from  $D$ ; see Lemma 23 and Definition 12;
- $\gamma_1 = \sqrt{2k}\gamma$ , which upper bounds  $|p_1(x)|$  for  $p_1 \in \mathcal{P}_{n, d, 2k}^1$  and  $x \in X_\gamma$ ; see Lemma 29;
- $\gamma_2 = 2\sqrt{s}\gamma^2$ , which upper bounds  $|p_2(x)|$  for  $p_2 \in \mathcal{P}_{n, d, s}^2$  and  $x \in X_\gamma$ ; see Lemma 30;
- $\rho_2 = C_2 \cdot d^{\frac{3}{4}} \cdot (c_0 d)^d$ , which upper bounds  $\|p_2\|_{L^2}$  for  $p_2 \in \mathcal{P}_{n, d, s}^2$ ; see Lemma 30;

## B. Omitted Proofs from Section 1 and Section 2

### B.1. Proof of Fact 3

*Proof.* It is known from Maass & Turán (1994) that in the absence of noise,  $\mathcal{H}_{d,K}$  can be PAC learned efficiently by using linear programming to find a concept that fits all the samples since in this case, empirical risk minimization with 0/1-loss is equivalent to solving a linear program. The number of samples is at least  $N_\delta := C_0 \cdot \frac{1}{\epsilon^2} (K^d \log n + \log \frac{1}{\delta})$  due to uniform convergence theory (Blumer et al., 1989) and the VC-dimension of  $\mathcal{H}_{d,K}$ . Then Theorem 12 of Kearns & Li (1988) shows that when  $\eta \leq \frac{1}{N_{1/2}} \log N_{1/2}$ , it is possible to learn the same concept class by using  $2N_\delta^2 \log \frac{1}{\delta} \cdot N_\delta = 2N_\delta^3 \log \frac{1}{\delta}$  samples. Substituting  $N_\delta$  gives the result.  $\square$

### B.2. Proof of Lemma 10

*Proof.* By definition, we know that there exists  $q \in \mathcal{P}_{n,d}$  and  $\Omega \subset [n]$  with  $|\Lambda| \leq K$ , such that  $f(x) = \text{sign}(q(x_\Lambda))$ . Let  $\bar{\Lambda} := [n] \setminus \Lambda$ . Since we choose Hermite polynomials as the basis, we have that  $m_i(x) = m_i(x_\Lambda) \cdot m_i(x_{\bar{\Lambda}})$  where we define  $m_i(x_{\bar{\Lambda}}) = 1$  if  $m_i(x)$  does not depend on  $x_{\bar{\Lambda}}$ .

We calculate the  $i$ -th element of the Chow vector of  $f$  as follows:

$$\mathbb{E}_{x \sim D}[f(x) \cdot m_i(x)] = \mathbb{E}_{x \sim D}[\text{sign}(q(x_\Lambda)) \cdot m_i(x)] \quad (\text{B.1})$$

$$= \mathbb{E}_{x \sim D}[\text{sign}(q(x_\Lambda)) \cdot m_i(x_\Lambda) \cdot m_i(x_{\bar{\Lambda}})] \quad (\text{B.2})$$

$$= \mathbb{E}_{x \sim D}[\text{sign}(q(x_\Lambda)) \cdot m_i(x_\Lambda)] \cdot \mathbb{E}_{x \sim D}[m_i(x_{\bar{\Lambda}})] \quad (\text{B.3})$$

as long as  $m_i$  depends on some elements in  $\bar{\Lambda}$ . Equivalently, the above is non-zero for all  $m_i$  that depends only on  $\Lambda$ . Note that there are at most

$$\sum_{i=0}^d K^i \leq \begin{cases} d+1, & \text{if } K=1, \\ \frac{K^{d+1}-1}{K-1} \leq 2K^d, & \text{if } K \geq 2 \end{cases} \quad (\text{B.4})$$

such  $m_i$ 's. This gives the desired sparsity bound.  $\square$

### B.3. The basis of monomials

As a complementary discussion to Lemma 10, we also give derivation for the sparsity of the Chow parameters under the basis of monomial polynomials.

**Lemma 22.** *Let  $f$  be a  $K$ -sparse degree- $d$  PTF. Then  $\chi_f$  is a  $k$ -sparse vector under the basis of monomial polynomials, where  $k \geq (\frac{n}{d})^{\lfloor \frac{d}{2} \rfloor}$ .*

*Proof.* Consider the same setting as that in Lemma 10, except that the basis is now under monomial polynomials. Since multivariate monomials are constructed by the production of univariate monomials, the  $i$ -th element of the Chow vector of  $f$  can be written as

$$\mathbb{E}_{x \sim D}[f(x) \cdot m_i(x)] = \mathbb{E}_{x \sim D}[\text{sign}(q(x_\Lambda)) \cdot m_i(x_\Lambda)] \cdot \mathbb{E}_{x \sim D}[m_i(x_{\bar{\Lambda}})].$$

However, now the term  $\mathbb{E}_{x \sim D}[m_i(x_{\bar{\Lambda}})]$  equals zero only when  $m_i(x_{\bar{\Lambda}})$  includes at least one univariate monomial  $x_j^\ell$  for some  $j \in \bar{\Lambda}$  where  $\ell \in \mathbb{Z}_+$  is an odd integer. For  $K \leq \frac{n}{2}$ ,  $d \geq 2$ , the total number of non-zero elements in  $\chi_f$  is at least

$$\sum_{j=1}^{\lfloor \frac{d}{2} \rfloor} \binom{n-K+j-1}{j} \geq \sum_{j=1}^{\lfloor \frac{d}{2} \rfloor} \left( \frac{n-K+j-1}{j} \right)^j \geq \left( \frac{n-K+\lfloor \frac{d}{2} \rfloor-1}{\lfloor \frac{d}{2} \rfloor} \right)^{\lfloor \frac{d}{2} \rfloor} \geq \left( \frac{n}{d} \right)^{\lfloor \frac{d}{2} \rfloor},$$

which depends polynomially in the dimension  $n$ , making it undesirable in the attribute-efficient learning.  $\square$

### B.4. Other choices of polynomial basis

As mentioned in Section 1.2, the choice of appropriate basis that demonstrates sparsity structure of the Chow parameters can go well beyond that of Hermite polynomials. More generally, under the assumption that  $D$  is a product distribution

(Eq.(B.2)), we require the basis to be a product basis (so Eq.(B.1) holds) with zero mean under the distribution  $D$  (so Eq.(B.2) holds). To this end, there exist other choices of basis under different distributional assumptions. For example, under the uniform distribution over  $[-1, 1]^n$ , the multivariate Legendre polynomials also form an appropriate basis. Another necessary property of  $D$  in our analysis is the finiteness of moments up to order  $4d$ , so that we can obtain tail bound for any degree- $4d$  polynomial; this is needed to establish Lemma 30.

## C. General Statistical Results

Recall that we assume  $D$  is the standard Gaussian  $\mathcal{N}(0, I_{n \times n})$  in this paper.

### C.1. Tail bound on $\|m(x)\|_\infty$

The tail bound of Fact 21 implies the following upper bound on the magnitude of  $m(x)$ .

**Lemma 23** (Restatement of Lemma 11). *The following holds for all  $t > 1$ :*

$$\begin{aligned} \Pr_{x \sim D} (\|m(x)\|_\infty \geq t) &\leq (n+1)^d \exp(-t^{2/d}/c_0), \\ \Pr_{S \sim D} (\max_{x \in S} \|m(x)\|_\infty \geq t) &\leq |S| \cdot (n+1)^d \cdot \exp(-t^{2/d}/c_0). \end{aligned}$$

In particular, with probability at least  $1 - \delta_\gamma$ , we have  $\max_{x \in S} \|m(x)\|_\infty \leq \gamma$  where  $\gamma := (c_0 \log \frac{|S|(n+1)^d}{\delta_\gamma})^{d/2}$ . When  $\delta_\gamma = \frac{\epsilon^2 \delta}{64 \rho_2^2}$  and  $|S| = C \cdot \frac{d^{5d} K^{4d}}{\epsilon^2} \log^{5d} \left( \frac{nd}{\epsilon \delta} \right)$ , we have  $\gamma = (C_1 d \cdot \log \frac{nd}{\epsilon \delta})^{d/2}$ .

*Proof.* Denote  $M = (n+1)^d$  the dimension of the vector  $m(x)$ . Let  $m_i(x)$  be the  $i$ -th element of  $m(x)$ , where  $1 \leq i \leq M$ . We note that  $m_1(x) \equiv 1$ . Now for any  $i \neq 1$ , since  $m(x)$  is orthonormal in  $L^2(\mathbb{R}^n, D)$ , we have  $\mathbb{E}[m_i(D)] = 0$  and  $\|m_i\|_{L^2} = 1$ . By Fact 21, we have

$$\Pr_{x \sim D} (|m_i(x)| \geq t) \leq \exp(-t^{2/d}/c_0).$$

Taking the union bound over all  $i \in \{2, \dots, M\}$  gives

$$\Pr_{x \sim D} \left( \max_{2 \leq i \leq M} |m_i(x)| \geq t \right) \leq (M-1) \exp(-t^{2/d}/c_0) \leq M \exp(-t^{2/d}/c_0).$$

Note that for all  $t > 1$ , we have

$$\Pr_{x \sim D} (\|m(x)\|_\infty \geq t) \leq M \exp(-t^{2/d}/c_0). \quad (\text{C.1})$$

Now for  $S$  being a set of independent draws from  $D$ , we have by union bound that

$$\Pr_{S \sim D} \left( \max_{x \in S} \|m(x)\|_\infty \geq t \right) \leq |S| \cdot M \cdot \exp(-t^{2/d}/c_0). \quad (\text{C.2})$$

The proof is complete.  $\square$

### C.2. $\beta(\tau, d, \rho), \beta'(\tau, d, \rho)$

**Lemma 24.** *Let  $Q_{d, \rho}(t) = \exp(-(t/\rho)^{2/d}/c_0)$  where  $\rho$  is independent of  $t$ . Then  $\int_0^\infty t \cdot \min\{\tau, Q_{d, \rho}(t)\} dt \leq \beta(\tau, d, \rho)$  where  $\beta(\tau, d, \rho) := \rho^2 \cdot (c_0 \log \frac{1}{\tau} + c_0 d)^d \cdot \tau$ .*

*Proof.* Denote  $\rho_0 = c_0 \cdot \rho^{2/d}$ . Let  $t_0 = (\rho_0 \log \frac{1}{\tau})^{d/2}$ , which satisfies  $\tau = Q_{d,\rho}(t_0)$ . It follows that

$$\begin{aligned}
 \int_0^\infty t \cdot \min\{\tau, Q_{d,\rho}(t)\} dt &= \int_0^{t_0} t \cdot \tau dt + \int_{t_0}^\infty t \cdot e^{-t^{2/d}/\rho_0} dt \\
 &\stackrel{\zeta_1}{=} \frac{1}{2} t_0^2 \tau + \frac{d \cdot \rho_0^d}{2} \int_{t_0^{2/d}/\rho_0}^\infty e^{-u} \cdot u^{d-1} du \\
 &\stackrel{\zeta_2}{\leq} \frac{1}{2} t_0^2 \tau + \frac{d \cdot \rho_0^d}{2} \cdot \Gamma(d, t_0^{2/d}/\rho_0) \\
 &\stackrel{\zeta_3}{\leq} \frac{1}{2} t_0^2 \tau + \frac{d \cdot \rho_0^d}{2} \cdot \exp(-t_0^{2/d}/\rho_0) \cdot (t_0^{2/d}/\rho_0 + d)^{d-1} \\
 &\stackrel{\zeta_4}{=} \frac{1}{2} \cdot \tau \cdot (\rho_0 \log \frac{1}{\tau})^d + \frac{d \cdot \rho_0^d}{2} \cdot \tau \cdot \left(\log \frac{1}{\tau} + d\right)^{d-1} \\
 &\leq \left(\rho_0 \log \frac{1}{\tau} + \rho_0 d\right)^d \cdot \tau.
 \end{aligned}$$

In the above,  $\zeta_1$  follows from the change of variables  $u := t^{2/d}/\rho$ ,  $\zeta_2$  follows from the definition of upper incomplete gamma function,  $\zeta_3$  follows from Lemma 28,  $\zeta_4$  follows from the setting of  $t_1$ , and the last step follows from simple algebraic relaxation. The result follows by replacing  $\rho_0$  with  $c_0 \cdot \rho^{2/d}$ .  $\square$

**Corollary 25.** *The following holds:*

$$\sup_{p_1 \in \mathcal{P}_{n,d,2k}^1} \int_0^\infty t \cdot \min\{\tau, \Pr(|p_1(D)| > t)\} dt \leq \beta(\tau, d, \sqrt{2}).$$

*Proof.* We will use the tail bound from Lemma 29. Let  $t_0 = \sqrt{2}(c_0 \log \frac{1}{\tau})^{d/2}$ . Since  $c_0 > 1$ , we have  $t_0 > \sqrt{2}$ . By Lemma 29, we also have  $\Pr(|p_1(D)| > t_0) \leq \tau$ . Therefore, we can write

$$\int_0^\infty t \cdot \min\{\tau, \Pr(p_1(D) \geq t)\} dt = \int_0^{t_0} t \cdot \tau dt + \int_{t_0}^\infty t \cdot e^{-(t/\sqrt{2})^{2/d}/c_0} dt.$$

Using the same proof as Lemma 24 gives the desired result.  $\square$

**Corollary 26.** *The following holds:*

$$\sup_{p_2 \in \mathcal{P}_{n,d,s}^2} \int_0^\infty t \cdot \min\{\tau, \Pr(|p_2(D)| > t)\} dt \leq \beta(\tau, 2d, \rho_2).$$

*Proof.* This follows immediately from the tail bound on  $\Pr(|p_2(D)| > t)$  in Lemma 30 and Lemma 24.  $\square$

The following lemma provides another type of bound.

**Lemma 27.** *Let  $Q_{d,\rho}(t) = \exp(-(t/\rho)^{2/d}/c_0)$  where  $\rho$  is independent of  $t$ . Then  $\int_0^\infty \min\{\tau, Q_{d,\rho}(t)\} dt \leq \beta'(\tau, d, \rho)$  where  $\beta'(\tau, d, \rho) := 2\rho \cdot (c_0 \log \frac{1}{\tau} + c_0 \cdot \frac{d}{2})^{d/2} \cdot \tau$ .*

*Proof.* Denote  $\rho_0 = c_0 \cdot \rho^{2/d}$ . Let  $t_0 = (\rho_0 \log \frac{1}{\tau})^{d/2}$ , which satisfies  $\tau = Q_{d,\rho}(t_0)$ . It follows that

$$\begin{aligned}
 \int_0^\infty t \cdot \min\{\tau, Q_{d,\rho}(t)\} dt &= \int_0^{t_0} \tau dt + \int_{t_0}^\infty e^{-t^{2/d}/\rho_0} dt \\
 &\stackrel{\zeta_1}{=} t_0 \tau + \rho_0^{d/2} \int_{t_0^{2/d}/\rho_0}^\infty e^{-u} \cdot u^{d/2-1} du \\
 &\stackrel{\zeta_2}{=} t_0 \tau + \rho_0^{d/2} \cdot \Gamma(d/2, t_0^{2/d}/\rho_0) \\
 &\stackrel{\zeta_3}{\leq} t_0 \tau + \rho_0^{d/2} \cdot \exp(-t_0^{2/d}/\rho_0) \cdot (t_0^{2/d}/\rho_0 + d/2)^{d/2-1} \\
 &\stackrel{\zeta_4}{=} \tau \cdot (\rho_0 \log \frac{1}{\tau})^{d/2} + \rho_0^{d/2} \cdot \tau \cdot \left(\log \frac{1}{\tau} + \frac{d}{2}\right)^{d/2-1} \\
 &\leq 2\rho_0^{d/2} \cdot \left(\log \frac{1}{\tau} + \frac{d}{2}\right)^{d/2} \cdot \tau.
 \end{aligned}$$

In the above,  $\zeta_1$  follows from the change of variables  $u := t^{2/d}/\rho$ ,  $\zeta_2$  follows from the definition of upper incomplete gamma function,  $\zeta_3$  follows from Lemma 28,  $\zeta_4$  follows from the setting of  $t_0$ , and the last step follows from simple algebraic relaxation. The result follows by noting  $\rho_0 = c_0 \cdot \rho^{2/d}$ .  $\square$

**Lemma 28** (Claim 3.11 of [Diakonikolas et al., 2018b]). *Consider the upper incomplete gamma function  $\Gamma(s, x) = \int_x^\infty t^{s-1} e^{-t} dt$ . We have  $\Gamma(s, x) \leq e^{-x} \cdot (x + s)^{s-1}$  for all  $s \geq 1$  and  $x \geq 0$ .*

### C.3. Concept class: basic properties

Recall that  $m(x)$  is the vector of all Hermite polynomials on  $\mathbb{R}^n$  with degree at most  $d$ . Note that  $m(x)$  has  $(n+1)^d$  elements. We defined two classes of polynomials:

$$\mathcal{P}_{n,d,2k}^1 := \{p : \mathbb{R}^n \rightarrow \mathbb{R} : p(x) = \langle v, m(x) \rangle, v \in \mathbb{R}^{(n+1)^d}, \|v\|_2 = 1, \|v\|_0 \leq k\}, \quad (\text{C.3})$$

and

$$\mathcal{P}_{n,d,s}^2 := \{p : \mathbb{R}^n \rightarrow \mathbb{R} : p = \langle A_U, mm^\top - I \rangle, A \in \mathbb{S}^{(n+1)^d}, U^\top = U, \|U\|_0 \leq s, \|A_U\|_F = 1\}, \quad (\text{C.4})$$

where

$$\mathbb{S}^{(n+1)^d} = \{A : A \in \mathbb{R}^{(n+1)^d \times (n+1)^d}, A^\top = A\}. \quad (\text{C.5})$$

Note that the polynomials in  $\mathcal{P}_{n,d,s}^2$  have degree at most  $2d$ , and can be represented as a linear combination of at most  $s$  elements of the form  $m_i(x)m_j(x)$ .

We collect a few basic properties of the sparse polynomial classes.

**Lemma 29.** *For all  $p_1 \in \mathcal{P}_{n,d,2k}^1$ , the following holds:*

- *Deterministic property:*  $|p_1(x)| \leq \sqrt{2k} \cdot \|m(x)\|_\infty$ ; in particular,  $|p_1(x)| \leq \gamma_1$  for all  $x \in X_\gamma$ , where  $\gamma_1 := \sqrt{2k}\gamma$ ;
- *Distributional property:* for all  $t \geq \sqrt{2}$ ,

$$\Pr(|p_1(D)| \geq t) \leq e^{-(t/\sqrt{2})^{2/d}/c_0}.$$

*Proof.* By Holder's inequality, we have

$$|p_1(x)| = |v \cdot m(x)| \leq \sqrt{\|v\|_0} \cdot \|v\|_2 \cdot \|m(x)\|_\infty \leq \sqrt{2k} \cdot 1 \cdot \|m(x)\|_\infty,$$

where the first inequality follows from (A.5) and the second inequality follows from the fact that  $p_1 \in \mathcal{P}_{n,d,2k}^1$ .

To show the tail bound, we will decompose  $v = (v_1, \tilde{v})$  and  $m(x) = (1, \tilde{m}(x))$  so that  $\mathbb{E}[\tilde{m}(D)] = 0$ . In this way, we have  $p_1(x) = v_1 + \tilde{v} \cdot \tilde{m}(x) = v_1 + \|\tilde{v}\|_2 \cdot q(x)$ , where  $q(x) := \tilde{v} \cdot \tilde{m}(x)$  and  $\bar{v} := \tilde{v} / \|\tilde{v}\|_2$ .

Observe that  $\mathbb{E}[q(D)] = 0$  and  $\text{Var}[q(D)] = 1$ . By Fact 21, for all  $t > 0$ ,

$$\Pr(|q(D)| \geq t) \leq e^{-t^{2/d}/c_0}. \quad (\text{C.6})$$

By simple calculation, we can show that

$$|p_1(x)| = \sqrt{(v_1 + \|\tilde{v}\|_2 \cdot q(x))^2} \leq \sqrt{2} \cdot \sqrt{v_1^2 + \|\tilde{v}\|_2^2 \cdot q^2(x)} \leq \sqrt{2} |q(x)|$$

for  $q(x) \geq 1$ . Thus, for all  $t \geq \sqrt{2}$ , we have

$$\Pr(|p_1(D)| \geq t) \leq \Pr(\sqrt{2} |q(D)| \geq t) = \Pr(|q(D)| \geq t/\sqrt{2}) \leq e^{-(t/\sqrt{2})^{2/d}/c_0}, \quad (\text{C.7})$$

where the last step follows from (C.6).  $\square$

**Lemma 30.** *For all  $p_2 \in \mathcal{P}_{n,d,s}^2$ , the following holds:*

- *Deterministic property:*  $|p_2(x)| \leq 2\sqrt{s} \cdot \|m(x)\|_\infty^2$ ; in particular,  $|p_2(x)| \leq \gamma_2$  for all  $x \in X_\gamma$ , where  $\gamma_2 := 2\sqrt{s} \cdot \gamma^2$ ;
- *Distributional properties:*  $\mathbb{E}[p_2(D)] = 0$ ,  $\|p_2\|_{L^2} \leq \rho_2$  where  $\rho_2 := C_2 \cdot d^{3/4} \cdot (c_0 d)^d$ , and

$$\Pr(|p_2(D)| \geq t) \leq \exp(-(t/\rho_2)^{1/d}/c_0), \quad \forall t > 0.$$

*Proof.* By the Cauchy-Schwartz inequality,

$$|p_2(x)| \leq \|A_U\|_F \cdot \left\| (m(x)m(x)^\top - I)_U \right\|_F \leq \sqrt{\|U\|_0} \cdot \|m(x)\|_\infty^2 + \|I_U\|_F \leq \sqrt{s} (\|m(x)\|_\infty^2 + 1),$$

where in the second inequality we use the fact that each entry of the matrix  $m(x)m(x)^\top$  takes the form  $m_i(x)m_j(x)$  whose magnitude is always upper bounded by  $\|m(x)\|_\infty^2$ , and there are at most  $\|U\|_0$  such entries. Since  $m_1(x) = 1$ , we always have  $1 \leq \|m(x)\|_\infty^2$ . Thus  $|p_2(x)| \leq 2\sqrt{s} \cdot \|m(x)\|_\infty^2$ .

Now we show the distributional properties. Since  $m(x)$  is a complete orthonormal basis in  $L^2(\mathbb{R}^n, D)$ , it follows that

$$\mathbb{E}[p_2(D)] = \langle A_U, I - I \rangle = 0. \quad (\text{C.8})$$

Now we bound  $\|p_2\|_{L^2}$ , which equals  $\sqrt{\mathbb{E}[p_2^2(D)]}$ . Recall that  $A_U$  is symmetric with  $\|A_U\|_F = 1$ , and thus can be written as

$$A_U = V^\top \Lambda V, \quad (\text{C.9})$$

for some orthonormal matrix  $V$  and diagonal matrix  $\Lambda$  with  $\|\Lambda\|_F = 1$ . Let  $q(x) = V \cdot m(x)$ . Observe that  $\mathbb{E}[q(D)] = 0$  and  $\mathbb{E}[q(D)q(D)^\top] = I$ . Then

$$\mathbb{E}[p_2^2(D)] = \text{Var}[q(D)^\top \Lambda q(D)] = \text{Var}\left[\sum_i \Lambda_{ii} q_i^2(D)\right] \leq \sum_i \Lambda_{ii}^2 \text{Var}[q_i^2(D)], \quad (\text{C.10})$$

where  $\Lambda_{ii}$  denotes the  $i$ -th diagonal element of  $\Lambda$  and  $q_i(x)$  denotes the  $i$ -th component of the vector-valued function  $q(x)$ , and the last step follows from the fact that  $\mathbb{E}[q_i(D) \cdot q_j(D)] = 0$  for  $i \neq j$  and  $\mathbb{E}[q_i(D)] = 0$  for all  $i$ . It thus remains to upper bound  $\text{Var}[q_i^2(D)]$ .

By the definition of variance, we have

$$\text{Var}[q_i^2(D)] = \mathbb{E}[q_i^4(D)] - (\mathbb{E}[q_i^2(D)])^2 = \mathbb{E}[q_i^4(D)] - 1 \leq C_2^2 \cdot d^{3/2} \cdot (c_0 d)^{2d},$$

where we applied Lemma 31 in the last step. Plugging it into (C.10), we get

$$\mathbb{E}[p_2^2(D)] \leq C_2^2 \cdot d^{3/2} \cdot \left(\frac{2d-1}{c_0 e}\right)^{2d-1} \sum_i \Lambda_{ii}^2 = C_2^2 \cdot d^{3/2} \cdot (c_0 d)^{2d}, \quad (\text{C.11})$$

where the equality follows from the construction that  $\|\Lambda\|_F = 1$  and  $\sum_i \Lambda_{ii}^2 = \|\Lambda\|_F^2$ . The proof is complete by noting that  $\|p_2\|_{L^2} = \sqrt{\mathbb{E}[p_2^2(D)]}$ .  $\square$

**Lemma 31.** Let  $v$  be a unit vector and  $m(x)$  be a collection of orthonormal polynomials of degree at most  $d$  in  $L^2(\mathbb{R}^n, D)$ . Then  $\mathbb{E}[(v \cdot m(D))^4] \leq C_2^2 \cdot d^{\frac{3}{2}} \cdot (c_0 d)^{2d}$  for some sufficiently large constant  $C_2 > 0$ .

*Proof.* Let  $p(x) = v \cdot m(x)$  and  $\rho = c_0 \cdot 2^{1/d}$ . We bound the desired expectation by using Fact 20 with the tail bound in Lemma 29.

$$\begin{aligned} \mathbb{E}[p^4(D)] &= \int_0^\infty \Pr(p^4(D) > t) dt \\ &= \int_0^{\sqrt{2}} \Pr(p^4(D) > t) dt + \int_{\sqrt{2}}^\infty \Pr(p^4(D) > t) dt \\ &\leq \sqrt{2} + 4 \int_0^\infty t^3 \cdot \Pr(|p(D)| > t) dt \\ &\leq \sqrt{2} + 4 \int_0^\infty t^3 \cdot e^{-t^{2/d}/\rho} dt \\ &= \sqrt{2} + \frac{d}{2} \cdot \rho^{2d-1} \cdot \int_0^\infty t^{2d-1} e^{-t} dt \\ &\stackrel{\zeta_1}{=} \sqrt{2} + \frac{d}{2} \cdot \rho^{2d-1} \cdot (2d-1)! \\ &\stackrel{\zeta_2}{\leq} \sqrt{2} + \frac{d}{2} \cdot \rho^{2d-1} \cdot \sqrt{2\pi(2d-1)} \left(\frac{2d-1}{e}\right)^{2d-1} \cdot e^{\frac{1}{12(2d-1)}} \end{aligned}$$

where  $\zeta_1$  follows from the known fact on the value of the Gamma function, and  $\zeta_2$  follows from the Stirling's approximation. The result follows by noting that  $\rho = c_0 \cdot 2^{1/d}$  and choosing a large enough constant  $C_2$ .  $\square$

#### C.4. Concept class: uniform convergence and sample complexity

**Lemma 32.** The VC-dimension of the class  $\mathcal{H}_{n,d,2k}^1 := \{h : x \mapsto \text{sign}(p_1(x)), p_1 \in \mathcal{P}_{n,d,2k}^1\}$  is  $O(d \cdot k \log n)$ , and that of the class  $\mathcal{H}_{n,d,s}^2 := \{h : x \mapsto \text{sign}(p_2(x)), p_2 \in \mathcal{P}_{n,d,s}^2\}$  is  $O(d \cdot s \log n)$ .

*Proof.* For the class  $\mathcal{H}_{n,d,k}^1$ , we can consider the class of polynomials in  $\mathcal{P}_{n,d,2k}^1$  with a fixed support set. It is easy to see that the VC-dimension of such class is  $k+1$ . Now note that the number of choices of such support set is

$$\sum_{i=0}^k \binom{(n+1)^d}{i} \leq \left(\frac{e(n+1)^d}{k}\right)^k.$$

The concept class union argument states that for any  $\mathcal{H} = \cup_{i=1}^M \mathcal{H}_i$ , the VC dimension of  $\mathcal{H}$  is upper bounded by  $O(\max\{V, \log M + V \log \frac{\log M}{V}\})$ , where  $V$  is an upper bound on the VC dimension of all  $\mathcal{H}_i$ . Thus, the VC-dimension of  $\mathcal{H}_{n,d,2k}^1$  is  $O(d \cdot k \log n)$  by algebraic calculations.

Likewise, for  $\mathcal{H}_{n,d,s}^2$ , we can first fix the support  $U \subset [(n+1)^d] \times [(n+1)^d]$  in the representation of  $p_2$ . Let  $\mathcal{P}(U)$  be the class of polynomials in  $\mathcal{P}_{n,d,s}^2$  with the fixed  $U$ . It is easy that the VC-dimension of  $\mathcal{P}(U)$  is  $s+1$ . Now note that the number of choices of  $U$  is

$$\sum_{i=0}^s \binom{(n+1)^{2d}}{i} \leq \left(\frac{e(n+1)^{2d}}{s}\right)^s.$$

Using the same argument gives that the VC-dimension of  $\mathcal{H}_{n,d,s}^2$  is  $O(d \cdot s \log n)$ .  $\square$

**Proposition 33** (Restatement of Prop. 13). Let  $S$  be a set of  $C \cdot \frac{d^{5d} K^{4d}}{\epsilon^2} \log^{5d} \left(\frac{nd}{\epsilon\delta}\right)$  instances drawn independently from  $D$ , where  $C > 0$  is a sufficiently large constant. Then with probability  $1 - \delta$ ,  $S$  is a good set in the sense of Definition 12.

*Proof.* By Lemma 11, our setting of  $\delta_\gamma$  and  $|S|$ , it follows that with probability at least  $1 - \delta_\gamma$ , we must have  $S|_{X_\gamma} = S$ . This proves Part 1. From now on, we condition on this event happening.

We note that by the classical VC theory (Anthony & Bartlett, 1999) and our VC-dimension upper bound in Lemma 32, Parts 2, 3, 5, 6 in Definition 12 all hold with probability  $1 - \delta/4$  as far as

$$|S| \geq C \cdot \left( \frac{\text{VCdim}}{\alpha^2} \cdot \log \frac{4 \cdot \text{VCdim}}{\alpha\delta} \right), \quad (\text{C.12})$$

where  $\text{VCdim} := \max\{d \cdot k \log n, d \cdot s \log n\} = d \cdot k^2 \log n$ , and  $\alpha := \min\{\alpha_1, \alpha_2\} \leq \frac{\epsilon}{2k\gamma^2}$ .

We now show Part 4. For  $x \in X_\gamma$ , we have  $|m_i(x)| \leq \gamma$  with certainty. Therefore, by Hoeffding's inequality for bounded random variables, we have that

$$\Pr \left( \left| \mathbb{E}_{S|X_\gamma} [f(x)m_i(x)] - \mathbb{E}_{D|X_\gamma} [f(x)m_i(x)] \right| > t \right) \leq 2 \exp \left( - \frac{|S|t^2}{4\gamma^2} \right).$$

Therefore, taking the union bound over  $i$  we obtain that if

$$|S| \geq \frac{32k\gamma^2}{(\alpha'_1)^2} \cdot \log \frac{16(n+1)^d}{\delta}, \quad (\text{C.13})$$

then with probability at least  $1 - \delta/8$ , we have

$$\max_{1 \leq i \leq (n+1)^d} \left| \mathbb{E}_{S|X_\gamma} [f(x)m_i(x)] - \mathbb{E}_{D|X_\gamma} [f(x)m_i(x)] \right| \leq \frac{1}{2} \cdot \frac{\alpha'_1}{\sqrt{2k}}.$$

Now we observe that for any  $p_1 \in \mathcal{P}_{n,d,2k}^1$ , we have  $p_1(x) = \langle v, m(x) \rangle$  with  $\|v\|_1 \leq \sqrt{k}$ . Thus

$$\begin{aligned} & \left| \mathbb{E}_{x \sim S|X_\gamma} [f(x) \cdot p_1(x)] - \mathbb{E}_{x \sim D|X_\gamma} [f(x) \cdot p_1(x)] \right| \\ &= \left| v \cdot \left( \mathbb{E}_{x \sim S|X_\gamma} [f(x) \cdot m(x)] - \mathbb{E}_{x \sim D|X_\gamma} [f(x) \cdot m(x)] \right) \right| \\ &\leq \sqrt{k} \cdot \left\| \mathbb{E}_{x \sim S|X_\gamma} [f(x) \cdot m(x)] - \mathbb{E}_{x \sim D|X_\gamma} [f(x) \cdot m(x)] \right\|_\infty \\ &\leq \frac{1}{2} \alpha'_1. \end{aligned} \quad (\text{C.14})$$

On the other hand, recall that for any  $p_1 \in \mathcal{P}_{n,d,2k}^1$ ,  $\|p_1\|_{L^2} = 1$  in view of Lemma 29. Thus Lemma 34 tells that

$$\sup_{f: |f| \leq 1, p_1 \in \mathcal{P}_{n,d,2k}^1} \left| \mathbb{E}_{x \sim D} [f(x) \cdot p(x)] - \mathbb{E}_{x \sim D|X_\gamma} [f(x) \cdot p(x)] \right| \leq 4 \sqrt{\Pr_{x \sim D} (x \notin X_\gamma)} \leq 4\sqrt{\delta_\gamma},$$

where the last step follows from Lemma 23. Since we set  $\delta_\gamma$  such that  $\delta_\gamma \leq \frac{(\alpha'_1)^2}{64}$ , we have

$$\sup_{f: |f| \leq 1, p_1 \in \mathcal{P}_{n,d,2k}^1} \left| \mathbb{E}_{x \sim D} [f(x) \cdot p(x)] - \mathbb{E}_{x \sim D|X_\gamma} [f(x) \cdot p(x)] \right| \leq \frac{1}{2} \alpha'_1. \quad (\text{C.15})$$

Part 4 follows from applying triangle inequality on (C.14) and (C.15), and the conditioning  $S|X_\gamma = S$ .

Lastly, we show Part 7. We note that for any fixed index  $(i, j) \in [(n+1)^d] \times [(n+1)^d]$ ,  $\sup_{x \in X_\gamma} |m_i(x)m_j(x)| \leq \gamma^2$  holds with certainty. Therefore, Hoeffding's inequality for bounded random variable tells that

$$\Pr \left( \left| \mathbb{E}_{S|X_\gamma} [m_i(x)m_j(x)] - \mathbb{E}_{D|X_\gamma} [m_i(x)m_j(x)] \right| > t \right) \leq 2 \exp \left( - \frac{|S|t^2}{4\gamma^4} \right).$$

Thus, by taking the union bound over all choices of  $(i, j)$ , we obtain that with probability  $1 - \delta/8$ ,

$$\max_{(i,j) \in [(n+1)^d] \times [(n+1)^d]} \left| \mathbb{E}_{S|X_\gamma} [m_i(x)m_j(x)] - \mathbb{E}_{D|X_\gamma} [m_i(x)m_j(x)] \right| \leq \frac{1}{2} \cdot \frac{\alpha'_2}{\sqrt{s}} \quad (\text{C.16})$$

as far as

$$|S| \geq \frac{16s\gamma^4}{(\alpha'_2)^2} \cdot \log \frac{16(n+1)^{2d}}{\delta}. \quad (\text{C.17})$$

Now, for any  $p_2 \in \mathcal{P}_{n,d,s}^2$ , we have  $p_2(x) = \langle A_U, m(x)m(x)^\top \rangle - \langle A_U, I \rangle$ . Thus,

$$\begin{aligned} & \left| \mathbb{E}_{S|X_\gamma} [p_2(x)] - \mathbb{E}_{D|X_\gamma} [p_2(x)] \right| \\ &= \left| \sum_{(i,j) \in U} A_{ij} (\mathbb{E}_{S|X_\gamma} [m_i(x)m_j(x)] - \mathbb{E}_{D|X_\gamma} [m_i(x)m_j(x)]) \right| \\ &\leq \sqrt{\|U\|_0} \cdot \max_{(i,j) \in [(n+1)^d] \times [(n+1)^d]} \left| \mathbb{E}_{S|X_\gamma} [m_i(x)m_j(x)] - \mathbb{E}_{D|X_\gamma} [m_i(x)m_j(x)] \right| \\ &\leq \sqrt{s} \cdot \frac{1}{2} \cdot \frac{\alpha'_2}{\sqrt{s}} \\ &= \frac{1}{2} \alpha'_2. \end{aligned} \quad (\text{C.18})$$

On the other hand, by Lemma 30, we have  $\|p_2\|_{L^2} \leq \rho_2$  for all  $p_2 \in \mathcal{P}_{n,d,s}^2$ . Thus, Lemma 34 tells that

$$\sup_{p_2 \in \mathcal{P}_{n,d,s}^2} \left| \mathbb{E}_{x \sim D} [p_2(x)] - \mathbb{E}_{x \sim D|X_\gamma} [p_2(x)] \right| \leq 4\rho_2 \sqrt{\Pr_{x \sim D} (x \notin X_\gamma)} \leq 4\rho_2 \sqrt{\delta_\gamma}.$$

Since we set  $\delta_\gamma$  such that  $\delta_\gamma \leq \frac{(\alpha'_2)^2}{64\rho_2^2}$ , we have

$$\sup_{p_2 \in \mathcal{P}_{n,d,s}^2} \left| \mathbb{E}_{x \sim D} [p_2(x)] - \mathbb{E}_{x \sim D|X_\gamma} [p_2(x)] \right| \leq 4\rho_2 \sqrt{\Pr_{x \sim D} (x \notin X_\gamma)} \leq \frac{1}{2} \alpha'_2. \quad (\text{C.19})$$

Part 7 follows from applying triangle inequality on (C.18) and (C.19), and the conditioning  $S|X_\gamma = S$ .

Observe that by the union bound, all these parts hold simultaneously with probability at least  $1 - \delta_\gamma - \delta/4 - \delta/8 - \delta/8 \geq 1 - \delta$  since we set  $\delta_\gamma \leq \delta/2$ . In addition, to satisfy all the requirements on the sample size involved in all parts, i.e. (C.12), (C.13), and (C.17), we need

$$|S| \geq C' \cdot \left( \frac{\gamma^4 \cdot k^4 \cdot \log n}{\epsilon^2} \log \frac{\gamma k n^d d}{\epsilon \delta} \right), \quad (\text{C.20})$$

for some large enough constant  $C' > 0$ . Our setting on  $|S|$  follows by plugging the setting of  $\gamma$  in Definition 12 into the above equation and noting that  $k \leq \max\{d+1, 2K^n\}$ . The proof is complete.  $\square$

**Lemma 34** (Total variation distance). *Assume that  $\Pr_{x \sim D} (x \in X_\gamma) \geq \frac{1}{2}$  and let  $\rho > 0$  be a finite real number. The following holds uniformly for all functions  $f$  and  $p$  satisfying  $f : \mathbb{R}^n \rightarrow [-1, 1]$  and  $\|p\|_{L^2} \leq \rho$ :*

$$\left| \mathbb{E}_{x \sim D} [f(x) \cdot p(x)] - \mathbb{E}_{x \sim D|X_\gamma} [f(x) \cdot p(x)] \right| \leq 4\rho \sqrt{\Pr_{x \sim D} (x \notin X_\gamma)}.$$

*Proof.* Denote  $z(x) = f(x) \cdot p(x)$ . Let  $\mathbf{1}_{X_\gamma^c}(x)$  be the indicator function which outputs 1 if  $x \notin X_\gamma$  and 0 otherwise. By simple calculation, we have

$$\mathbb{E}_{x \sim D} [z(x)] = \Pr_{x \sim D} (x \in X_\gamma) \cdot \mathbb{E}_{D|X_\gamma} [z(x)] + \mathbb{E}_D [z(x) \cdot \mathbf{1}_{X_\gamma^c}(x)],$$

namely,

$$\mathbb{E}_{D|X_\gamma} [z(x)] = \frac{\mathbb{E}_D [z(x)]}{\Pr_D (x \in X_\gamma)} - \frac{\mathbb{E}_D [z(x) \cdot \mathbf{1}_{X_\gamma^c}(x)]}{\Pr_D (x \in X_\gamma)}. \quad (\text{C.21})$$

Therefore,

$$\begin{aligned}
 \left| \mathbb{E}_{D|X_\gamma}[z(x)] - \mathbb{E}_D[z(x)] \right| &= \left| \frac{\Pr_D(x \notin X_\gamma) \cdot \mathbb{E}_D[z(x)]}{\Pr_D(x \in X_\gamma)} - \frac{\mathbb{E}_D[z(x) \cdot \mathbf{1}_{X_\gamma^c}(x)]}{\Pr_D(x \in X_\gamma)} \right| \\
 &\leq \frac{\Pr_D(x \notin X_\gamma)}{\Pr_D(x \in X_\gamma)} \cdot \left| \mathbb{E}_D[z(x)] \right| + \frac{\left| \mathbb{E}_D[z(x) \cdot \mathbf{1}_{X_\gamma^c}(x)] \right|}{\Pr_D(x \in X_\gamma)} \\
 &\leq \frac{1}{1/2} \cdot \Pr_D(x \notin X_\gamma) \cdot \left| \mathbb{E}_D[z(x)] \right| + \frac{\left| \mathbb{E}_D[z(x) \cdot \mathbf{1}_{X_\gamma^c}(x)] \right|}{1/2},
 \end{aligned} \tag{C.22}$$

where we applied the condition  $\Pr_D(x \in X_\gamma) \geq 1/2$  in the last step.

Observe that

$$\left| \mathbb{E}_D[z(x)] \right| \leq \sqrt{\mathbb{E}_D[z^2(x)]} \leq \sqrt{\mathbb{E}_D[p^2(x)]} = \|p\|_{L^2} \leq \rho. \tag{C.23}$$

On the other hand,

$$\left| \mathbb{E}_D[z(x) \cdot \mathbf{1}_{X_\gamma^c}(x)] \right| \leq \sqrt{\mathbb{E}_D[z^2(x)]} \cdot \sqrt{\mathbb{E}_D[\mathbf{1}_{X_\gamma^c}(x)]} \leq \rho \cdot \sqrt{\Pr_D(x \notin X_\gamma)}. \tag{C.24}$$

Combining (C.22), (C.23), (C.24), and noting that any probability is always no greater than its root completes the proof.  $\square$

## D. Analysis of Algorithm 2: Proof of Theorem 14

*Proof of Theorem 14.* Note that in view of our construction of  $p_2$  in the algorithm, we have  $\mathbb{E}[p_2(S')] = \|(\Sigma - I)_U\|_F$ . Denote  $E = S' \setminus S$  and  $L = S \setminus S'$ . Then,

$$|S'| \cdot \|(\Sigma - I)_U\|_F = |S'| \cdot \mathbb{E}[p_2(S')] = |S| \cdot \mathbb{E}[p_2(S)] + |E| \cdot \mathbb{E}[p_2(E)] - |L| \cdot \mathbb{E}[p_2(L)]. \tag{D.1}$$

Observe that Lemma 30 tells that  $\mathbb{E}[p_2(D)] = 0$ , which combined with Part 7 of Definition 12 gives  $\mathbb{E}[p_2(S)] \leq \alpha'_2$ . In addition, Lemma 36 shows  $|L| \cdot |\mathbb{E}[p_2(L)]| \leq 2(1 + \frac{1}{c})|S| \cdot (\beta'(\eta, 2d, \rho_2) + \alpha_2\gamma_2)$ . Assume for contradiction that no such threshold  $t$  exists. Then Lemma 35 gives  $|E| \cdot |\mathbb{E}[p_2(E)]| \leq 7(1 + \frac{1}{c})|S'| \cdot (\beta'(\eta, 2d, \rho_2) + \alpha_2\gamma_2)$ . Plugging these into (D.1), we obtain that

$$|S'| \cdot \|(\Sigma - I)_U\|_F \leq |S| \cdot \alpha'_2 + 7(1 + \frac{1}{c})|S'| \cdot (\beta'(\eta, 2d, \rho_2) + \alpha_2\gamma_2) + 2(1 + \frac{1}{c})|S| \cdot (\beta'(\eta, 2d, \rho_2) + \alpha_2\gamma_2).$$

Diving both sides by  $|S'|$  and noting that (A.2) shows  $|S| \leq (1 + \frac{1}{2c})|S'|$ , we obtain

$$\|(\Sigma - I)_U\|_F \leq 7(1 + \frac{1}{c})(1 + \frac{1}{2c})(\beta'(\eta, 2d, \rho_2) + \alpha'_2 + \alpha_2\gamma_2) \leq \frac{14}{c^2}(\beta'(\eta, 2d, \rho_2) + \alpha'_2 + \alpha_2\gamma_2),$$

where the last step follows since  $c \in (0, \frac{1}{2}]$ . Recall that we set  $\alpha'_2 = \epsilon$  and  $\alpha_2 = \epsilon/\gamma_2$  in Definition 12. Thus, the above inequality reads as

$$\|(\Sigma - I)_U\|_F \leq \frac{14}{c^2}(\beta'(\eta, 2d, \rho_2) + 2\epsilon) = \kappa,$$

which contradicts the condition of the proposition that  $\|(\Sigma - I)_U\|_F > \kappa$ .

Note that the existence of such threshold  $t$  combined with Lemma 38 implies the desired progress in the symmetric difference. In particular, by combining Part 5 of Definition 12 and Lemma 30, we have

$$\Pr(p_2(S) > t) \leq \exp(-(t/\rho_2)^{1/d}/c_0) + \alpha_2, \quad \forall t > 0. \tag{D.2}$$

We also just showed that there exists  $t > 0$  such that

$$\Pr(p_2(S') > t) \geq 6 \exp(-(t/\rho_2)^{1/d}/c_0) + 6\alpha_2. \tag{D.3}$$

In addition, (A.2) tells  $|S'| \geq \frac{1}{2}|S|$ . Thus, Lemma 38 asserts that

$$\Delta(S, S'') \leq \Delta(S, S') - \exp(-(t/\rho_2)^{1/d}/c_0) - \alpha_2 \leq \Delta(S, S') - \alpha_2. \tag{D.4}$$

This completes the proof by noting that we set  $\alpha_2 = \frac{\epsilon}{\gamma_2} = \frac{\epsilon}{4k\gamma^2}$ .  $\square$

### D.1. Auxiliary results

**Lemma 35** (Restatement of Lemma 15). *Consider Algorithm 2. Suppose that  $\Delta(S, S') \leq 2\eta$  and  $\|(\Sigma - I)_U\|_F > \kappa$ . Let  $E = S' \setminus S$ . If there does not exist a threshold  $t > 0$  that satisfies Step 3, then*

$$(|E| / |S'|) \sup_{p_2 \in \mathcal{P}_{n,d,s}^2} \mathbb{E}[|p_2(E)|] \leq 7\left(1 + \frac{1}{c}\right) \cdot [\beta'(\eta, 2d, \rho_2) + \alpha_2 \gamma_2].$$

*Proof.* We use Lemma 37 to establish the result. We note that  $\Delta(S, S') \leq 2\eta$  implies  $|E| \leq \frac{\eta}{c} |S'|$  by (A.3). Since we assumed that no threshold  $t$  satisfies the filtering condition, we have

$$\Pr(|p_2(S')| > t) \leq 6 \exp(-(t/\rho_2)^{1/d}/c_0) + 6\alpha_2, \quad \forall t > 0.$$

By Lemma 27, we have  $\int_0^\infty \min\{\eta, \exp(-(t/\rho_2)^{1/d}/c_0)\} \leq \beta'(\eta, 2d, \rho_2)$ . Lastly, by Lemma 30, we have  $\max_{x \in S'} |p_2(x)| \leq \gamma_2$ . Thus, using Lemma 37 gives the result.  $\square$

**Lemma 36** (Restatement of Lemma 16). *Consider Algorithm 2. Suppose that  $S$  is a good set and  $\Delta(S, S') \leq 2\eta$ . We have*

$$(|L| / |S|) \sup_{p_2 \in \mathcal{P}_{n,d,s}^2} \mathbb{E}[|p_2(L)|] \leq 2\left(1 + \frac{1}{c}\right) [\beta'(\eta, 2d, \rho_2) + \alpha_2 \gamma_2].$$

*Proof.* We use Lemma 37 to establish the result. Similar to Lemma 35, we can show that  $|L| \leq \frac{\eta}{c} |S|$  by (A.3). By combining Lemma 30 and Part 5 of Definition 12, we have

$$\Pr(|p_2(S)| > t) \leq \exp(-(t/\rho_2)^{1/d}/c_0) + \alpha_2, \quad \forall t > 0.$$

By Lemma 27, we have  $\int_0^\infty \min\{\eta, \exp(-(t/\rho_2)^{1/d}/c_0)\} \leq \beta'(\eta, 2d, \rho_2)$ . Lastly, by Lemma 30, we have  $\max_{x \in S'} |p_2(x)| \leq \gamma_2$ . Thus, using Lemma 37 gives the result.  $\square$

The following lemma borrows from Lemma 2.10 of Diakonikolas et al. (2018a); the proof is included for completeness.

**Lemma 37.** *Let  $c_0 > 0$  be an absolute constant. Let  $S_0$  be a set of instances in  $\mathbb{R}^n$  and  $S_1 \subset S_0$ , with  $|S_1| \leq \omega_1 \tau |S_0|$  for some  $\omega_1, \tau > 0$ . Let  $p$  be such that  $\Pr(|p(S_0)| > t) \leq \omega_2 \cdot Q_d(t) + \alpha_0$  for all  $t \geq t_0$ , where  $\omega_2, Q_d(t), \alpha_0 > 0$ . Assume  $\max_{x \in S_0} |p(x)| \leq \gamma_0$ . Further assume that  $\int_0^\infty \min\{\tau, Q_d(t)\} dt \leq \beta_0$ . Then*

$$(|S_1| / |S_0|) \cdot \mathbb{E}[|p(S_1)|] \leq \omega_1 t_0 \cdot \tau + (\omega_1 + 1)(\omega_2 + 1)\beta_0 + \alpha_0 \gamma_0.$$

*Proof.* Since  $S_1 \subset S_0$ , we have  $|S_1| \cdot \Pr(|p(S_1)| > t) \leq |S_0| \cdot \Pr(|p(S_0)| > t)$ . Thus,

$$\Pr(|p(S_1)| > t) \leq \min\left\{1, \frac{|S_0|}{|S_1|} \cdot \Pr(|p(S_0)| > t)\right\}. \quad (\text{D.5})$$

By Fact 20, we have

$$\begin{aligned}
 (|S_1| / |S_0|) \cdot \mathbb{E}[|p(S_1)|] &\leq \int_0^\infty (|S_1| / |S_0|) \Pr(|p(S_1)| > t) dt \\
 &\stackrel{\zeta_1}{=} \int_0^{\gamma_0} (|S_1| / |S_0|) \Pr(|p(S_1)| > t) dt \\
 &\stackrel{\zeta_2}{\leq} \int_0^{\gamma_0} \min \left\{ |S_1| / |S_0|, \Pr(|p(S_0)| > t) \right\} dt \\
 &\stackrel{\zeta_3}{\leq} \int_0^{\gamma_0} \min \left\{ \omega_1 \tau, \Pr(|p(S_0)| > t) \right\} dt \\
 &\stackrel{\zeta_4}{\leq} \int_0^{t_0} \min \{ \omega_1 \tau, 1 \} dt + \int_{t_0}^{\gamma_0} \min \{ \omega_1 \tau, \omega_2 \cdot Q_d(t) + \alpha_0 \} dt \\
 &\stackrel{\zeta_5}{\leq} \omega_1 \tau t_0 + \int_{t_0}^{\gamma_0} \min \{ \omega_1 \tau, \omega_2 \cdot Q_d(t) \} dt + \int_{t_0}^{\gamma_0} \alpha_0 dt \\
 &\stackrel{\zeta_6}{\leq} \omega_1 t_0 \cdot \tau + (\omega_1 + 1)(\omega_2 + 1) \int_{t_0}^{\gamma_0} \min \{ \tau, Q_d(t) \} dt + \alpha_0(\gamma_0 - t_0) \\
 &\stackrel{\zeta_7}{\leq} \omega_1 t_0 \cdot \tau + (\omega_1 + 1)(\omega_2 + 1) \beta_0 + \alpha_0 \gamma_0.
 \end{aligned}$$

In the above,  $\zeta_1$  follows from the condition that  $p(x) \leq \gamma_0$  for all  $x \in S_1$ ,  $\zeta_2$  follows from (D.5),  $\zeta_3$  uses the condition  $|S_1| \leq \omega_1 \tau |S_0|$ ,  $\zeta_4$  uses the condition of the tail bound of  $p(S_0)$  when  $t \geq t_0$ ,  $\zeta_5$  applies elementary facts that  $\min\{\omega_1 \tau, 1\} \leq \omega_1 \tau$  and  $\min\{a, b + c\} \leq \min\{a, b\} + c$  for any  $c > 0$ ,  $\zeta_6$  uses the fact that both  $\frac{\omega_1}{(\omega_1+1)(\omega_2+1)}$  and  $\frac{\omega_1}{(\omega_1+1)(\omega_2+1)}$  are less than 1 for positive  $\omega_1$  and  $\omega_2$ , and  $\zeta_7$  applies the condition on the integral and uses the fact that both  $\tau$  and  $Q_d(t)$  are positive.  $\square$

The following lemma is implicit in prior works but we give a slightly more general statement; see e.g. Claim 5.13 of Diakonikolas et al. (2016).

**Lemma 38.** *Let  $S$  and  $S'$  be two instance sets with  $|S'| \geq \alpha |S|$  for some  $\alpha \in (0, 1]$ . Suppose that there exists  $t_0 > 0$  such that  $\Pr(g(S) \geq t_0) \leq h_1(t_0)$ ,  $\Pr(g(S') \geq t_0) > h_2(t_0)$ , and  $h_2(t_0) \geq \frac{3}{\alpha} \cdot h_1(t_0)$ . Let  $S'' = S' \cap \{x : g(x) \geq t_0\}$ . Then  $\Delta(S, S'') - \Delta(S, S') \leq -h_1(t_0)$ .*

*Proof.* Write  $E := S' \setminus S$  and  $L := S \setminus S'$ . Then  $S' = S \cup E \setminus L$ . Likewise, write  $E' := S'' \setminus S$  and  $L' := S \setminus S''$ . Then  $S'' = S \cup E' \setminus L'$ . Since  $S'' \subset S'$ , we have  $E' \subset E$  and  $L' \supset L$ . It is not hard to see that

$$\Delta(S, S'') - \Delta(S, S') = \frac{|E'| + |L'|}{|S|} - \frac{|E| + |L|}{|S|} = \frac{1}{|S|} \cdot (|L' \setminus L| - |E \setminus E'|). \quad (\text{D.6})$$

Let  $V := \{x : g(x) \geq t_0\}$ . By our assumption, it follows that

$$|S \cap V| \leq h_1(t_0) \cdot |S|, \quad |S' \cap V| > h_2(t_0) \cdot |S'|.$$

By basic set operations, we have  $E \setminus E' = (S' \setminus S) \cap V = (S' \cap V) \setminus S = (S' \cap V) \setminus (S \cap V)$ . Thus,

$$|E \setminus E'| \geq |S' \cap V| - |S \cap V| \geq h_2(t_0) \cdot |S'| - h_1(t_0) |S| \geq (\alpha \cdot h_2(t_0) - h_1(t_0)) |S|. \quad (\text{D.7})$$

On the other hand,  $L' \setminus L = (S' \cap S) \cap V$ . Thus,

$$|L' \setminus L| \leq |S \cap V| \leq h_1(t_0) \cdot |S|. \quad (\text{D.8})$$

Combining (D.7) and (D.8), and the condition of  $h_2(t_0) \geq \frac{3}{\alpha} \cdot h_1(t_0)$ , we have

$$|E \setminus E'| \geq 2h_1(t_0) \cdot |S| \geq |L' \setminus L| + h_1(t_0) \cdot |S|.$$

This combined with (D.6) completes the proof.  $\square$

## E. Performance Guarantees on the Output of Algorithm 1

### E.1. Proof of Theorem 17

*Proof of Theorem 17* We first show the following holds:

$$\sup_{p_1 \in \mathcal{P}_{n,d,2k}^1} \left| \mathbb{E}_{(x,y) \sim \bar{S}_l'} [y \cdot p_1(x)] - \mathbb{E}_{x \sim D} [f^*(x) \cdot p_1(x)] \right| \leq \frac{64}{c^2} \sqrt{\eta(\beta_\eta + \beta_\epsilon)} + \frac{\epsilon}{2}. \quad (\text{E.1})$$

To ease notation, write  $S' := S'_l$ ,  $L = S \setminus S'$ ,  $E = S' \setminus S$ . Let  $p_1$  be an arbitrary polynomial in  $\mathcal{P}_{n,d,2k}^1$ . As  $S' = S \cup E \setminus L$ , it is easy to see that

$$\begin{aligned} & |S'| \cdot \left| \mathbb{E}_{(x,y) \sim \bar{S}'} [y \cdot p_1(x)] - \mathbb{E}_{(x,y) \sim \bar{S}} [y \cdot p_1(x)] \right| \\ &= \left| (|S'| - |S|) \mathbb{E}_{\bar{S}} [y \cdot p_1(x)] + |L| \cdot \mathbb{E}_{\bar{L}} [y \cdot p_1(x)] - |E| \cdot \mathbb{E}_{\bar{E}} [y \cdot p_1(x)] \right| \\ &\leq \left| |S'| - |S| \right| \cdot \left| \mathbb{E}_{\bar{S}} [y \cdot p_1(x)] \right| + |L| \cdot \left| \mathbb{E}_{\bar{L}} [y \cdot p_1(x)] \right| + |E| \cdot \left| \mathbb{E}_{\bar{E}} [y \cdot p_1(x)] \right|. \end{aligned}$$

Note that the Cauchy–Schwarz inequality states that  $\mathbb{E}[y \cdot p_1(x)] \leq \sqrt{\mathbb{E}[y^2]} \cdot \sqrt{\mathbb{E}[p_1^2(x)]} = \sqrt{\mathbb{E}[p_1^2(x)]}$  where the last step follows since  $y \in \{-1, 1\}$ . Therefore, continuing the above inequality, we have

$$\begin{aligned} & |S'| \cdot \left| \mathbb{E}_{(x,y) \sim \bar{S}'} [y \cdot p_1(x)] - \mathbb{E}_{(x,y) \sim \bar{S}} [y \cdot p_1(x)] \right| \\ &\leq \left| |S'| - |S| \right| \cdot \sqrt{\mathbb{E}[p_1^2(S)]} + |L| \cdot \sqrt{\mathbb{E}[p_1^2(L)]} + |E| \cdot \sqrt{\mathbb{E}[p_1^2(E)]} \\ &\leq \left| |S'| - |S| \right| \cdot \sqrt{1 + 2\beta_{\delta_\gamma} + \epsilon} + \sqrt{|L| \cdot |S| \cdot (12\beta_\eta + 4\eta + \epsilon)} + \sqrt{\frac{6}{c} |E| \cdot |S'| (\kappa + \beta_\eta + \beta_{\delta_\gamma} + \eta + \epsilon)} \quad (\text{E.2}) \end{aligned}$$

where in the last step we applied Lemma 40, Lemma 41, Lemma 42, and denoted  $\beta_{\delta_\gamma} = \beta(2\delta_\gamma, d, \sqrt{2})$  and  $\beta_\eta = \beta(\eta, d, \sqrt{2})$ .

On the other hand, (A.2) implies

$$\left| |S'| - |S| \right| \leq \frac{2\eta}{1 - 2\eta} |S'| \leq \frac{\eta}{c} |S'|$$

for  $\eta \in [0, \frac{1}{2} - c]$ . We also have the following estimates:  $\max\{|E|, |L|\} \leq \eta |S| \leq \frac{\eta}{1 - 2\eta} |S'| \leq \frac{\eta}{2c} |S'|$ . Plugging these into (E.2), we have

$$\left| \mathbb{E}_{(x,y) \sim \bar{S}'} [y \cdot p_1(x)] - \mathbb{E}_{(x,y) \sim \bar{S}} [y \cdot p_1(x)] \right| \leq \frac{1}{c} \left[ \eta \sqrt{1 + \beta_{\delta_\gamma} + \epsilon} + 4 \sqrt{\eta(\kappa + \beta_{\delta_\gamma} + \beta_\eta + \eta + \epsilon)} \right]. \quad (\text{E.3})$$

On the other hand, we note that in view of Part 4 of Definition 12, we have

$$\left| \mathbb{E}_{(x,y) \sim \bar{S}} [y \cdot p_1(x)] - \mathbb{E}_{x \sim D} [f^*(x)p_1(x)] \right| = \left| \mathbb{E}_{x \sim S} [f^*(x)p_1(x)] - \mathbb{E}_{x \sim D} [f^*(x)p_1(x)] \right| \leq \alpha'_1, \quad (\text{E.4})$$

where the first step follows from the condition that  $f^*(\cdot)$  is the underlying PTF and  $\bar{S}$  is an uncorrupted sample set (which implies  $y = f^*(x)$  for any  $(x, y) \in \bar{S}$ ). By applying triangle inequality on (E.3) and (E.4), we have

$$\left| \mathbb{E}_{(x,y) \sim \bar{S}'} [y \cdot p_1(x)] - \mathbb{E}_{x \sim D} [f^*(x) \cdot p_1(x)] \right| \leq \frac{4}{c} \left[ \eta \sqrt{1 + \beta_{\delta_\gamma} + \epsilon} + \sqrt{\eta(\kappa + \beta_{\delta_\gamma} + \beta_\eta + \eta + \epsilon)} \right] + \alpha'_1. \quad (\text{E.5})$$

Now recall that  $\alpha'_1 = \epsilon/6$ ,  $\delta_\gamma$  is such that  $\beta_{\delta_\gamma} \leq \beta_\epsilon$ ,  $\eta \leq \beta_\eta$ , and  $\kappa \leq \frac{14}{c^2}(\beta_\eta + \epsilon)$ . Thus, by rearrangement, we have

$$\begin{aligned} & \left| \mathbb{E}_{(x,y) \sim \bar{S}'} [y \cdot p_1(x)] - \mathbb{E}_{x \sim D} [f^*(x) \cdot p_1(x)] \right| \\ &\leq \frac{16}{c^2} \left[ \eta \sqrt{1 + \beta_\epsilon + \epsilon} + \sqrt{\eta(\beta_\eta + \beta_\epsilon + \epsilon)} \right] + \frac{\epsilon}{6} \\ &\leq \frac{\zeta_1 32}{c^2} \sqrt{\eta(\eta + \eta\beta_\epsilon + \eta\epsilon + \beta_\eta + \beta_\epsilon + \epsilon)} + \frac{\epsilon}{6} \\ &\leq \frac{\zeta_2 64}{c^2} \sqrt{\eta(\beta_\eta + \beta_\epsilon)} + \frac{\epsilon}{6}, \end{aligned}$$

where in  $\zeta_1$  we used the elementary inequality  $\sqrt{a} + \sqrt{b} \leq 2\sqrt{a+b}$ , and in  $\zeta_2$  we used the fact that  $\eta \leq \beta_\eta$ ,  $\eta\epsilon < \epsilon \leq \beta_\epsilon$ . This proves (E.1) since the above holds for any  $p_1 \in \mathcal{P}_{n,d,2k}^1$ .

Now we note that for any  $p_1 \in \mathcal{P}_{n,d,2k}^1$ , it can be represented as  $p_1(x) = \langle v, m(x) \rangle$  with  $\|v\|_2 = 1$  and  $\|v\|_0 \leq 2k$ . In this way, we get

$$\begin{aligned} & \left| \mathbb{E}_{(x,y) \sim \bar{S}'}[y \cdot p_1(x)] - \mathbb{E}_{x \sim D}[f^*(x) \cdot p_1(x)] \right| \\ &= \left| \mathbb{E}_{(x,y) \sim \bar{S}'}[y \cdot \langle v, m(x) \rangle] - \mathbb{E}_{x \sim D}[f^*(x) \cdot \langle v, m(x) \rangle] \right| \\ &= \left| \langle v, \mathbb{E}_{(x,y) \sim \bar{S}'}[y \cdot m(x)] \rangle - \langle v, \mathbb{E}_{x \sim D}[f^*(x) \cdot m(x)] \rangle \right| \\ &= \left| \langle v, \mathbb{E}_{(x,y) \sim \bar{S}'}[y \cdot m(x)] \rangle - \chi_{f^*} \right|. \end{aligned}$$

Using Lemma 45 completes the proof.  $\square$

## E.2. Proof of Theorem 18

*Proof of Theorem 18.* Let  $\bar{S}$  be the uncorrupted sample set with the same size as  $\bar{S}'$ . Observe that by Proposition 13,  $S$  is a good set and  $\Delta(S, S') \leq 2\eta$ . We show by induction the progress of filtering, which will imply that within  $l_{\max}$  phases, Algorithm I must terminate.

Suppose that the algorithm returns at some phase  $\bar{l} \geq 1$ , i.e.  $\|(\Sigma - I)_U\|_F > \kappa$  for all  $1 \leq l < \bar{l}$ . We show by induction that the two invariants hold:  $\Delta(S, S'_l) \leq 2\eta$  and  $\Delta(S, S'_{l+1}) \leq \Delta(S, S'_l) - \frac{\epsilon}{2k\gamma^2}$ .

**Base case:**  $l = 1$ . Note that since  $S$  is a good set,  $S|_{X_\gamma} = S$ . Thus, no samples in  $S$  are pruned in Step I of Algorithm I. Therefore, we have  $\Delta(S, S'_1) \leq \Delta(S, S') \leq 2\eta$ . In addition, we have  $\|(\Sigma - I)_U\|_F > \kappa$ . Thus, Theorem 14 tells us that

$$\Delta(S, S'_2) \leq \Delta(S, S'_1) - \frac{\epsilon}{2k\gamma^2}. \quad (\text{E.6})$$

In particular, the above implies that  $\Delta(S, S'_2) \leq 2\eta$ .

**Induction.** Now suppose that  $\Delta(S, S'_l) \leq 2\eta$ . Then applying Theorem 14 gives us

$$\Delta(S, S'_{l+1}) \leq \Delta(S, S'_l) - \frac{\epsilon}{2k\gamma^2}, \quad (\text{E.7})$$

and in particular,  $\Delta(S, S'_{l+1}) \leq 2\eta$ .

Therefore, by telescoping, we obtain that

$$\Delta(S, S_{\bar{l}}) \leq \Delta(S, S'_1) - \frac{(\bar{l}-1) \cdot \epsilon}{2k\gamma^2} \leq 2\eta - \frac{(\bar{l}-1) \cdot \epsilon}{2k\gamma^2}. \quad (\text{E.8})$$

On the other hand, the symmetric difference  $\Delta(S, S_{\bar{l}})$  is always non-negative. This implies that

$$\bar{l} \leq \frac{4\eta k \gamma^2}{\epsilon} + 1 = \frac{4\eta k}{\epsilon} \cdot \left( C_1 d \cdot \log \frac{nd}{\epsilon \delta} \right)^d + 1, \quad (\text{E.9})$$

where we realized the setting of  $\gamma$  in the second step.

Now we characterize the output of the algorithm. In fact, by Theorem 17, we have

$$\left\| H_k \left( \mathbb{E}_{(x,y) \sim \bar{S}'_l} [y \cdot m(x)] \right) - \chi_{f^*} \right\|_2 \leq \frac{192}{c^2} \sqrt{\eta(\beta_\eta + \beta_\epsilon)} + \frac{\epsilon}{2}.$$

The proof is complete by noting that  $H_k \left( \mathbb{E}_{(x,y) \sim \bar{S}'_l} [y \cdot m(x)] \right)$  is the output of the algorithm.  $\square$

### E.3. Auxiliary results

**Lemma 39.** *If  $\|(\Sigma - I)_U\|_F \leq \kappa$ , then*

$$\sup_{p_1 \in \mathcal{P}_{n,d,2k}^1} \mathbb{E}[p_1^2(S')] \leq \kappa + 1.$$

*Proof.* For any  $p_1 \in \mathcal{P}_{n,d,2k}$ , we can write it as  $p_1(x) = v \cdot m(x)$  where  $v$  is  $2k$ -sparse and  $\|v\|_2 = 1$ . Denote by  $J$  the support set of  $v$ . Then,

$$\begin{aligned} \mathbb{E}[p_1^2(S')] &= \mathbb{E}[v^\top m(S') m^\top(S') v] \\ &= v^\top \Sigma_{J \times J} v = v^\top (\Sigma - I)_{J \times J} v + v^\top v \leq \|vv^\top\|_F \cdot \|(\Sigma - I)_{J \times J}\|_F + 1. \end{aligned}$$

Since  $\|v\|_0 \leq 2k$ , we know that  $J \times J$  has  $2k$  diagonal entries and  $4k^2 - 2k$  off-diagonal symmetric entries. This implies  $\|(\Sigma - I)_{J \times J}\|_F \leq \|(\Sigma - I)_U\|_F \leq \kappa$ . Now using  $\|vv^\top\|_F = \|v\|_2^2 = 1$  completes the proof.  $\square$

**Lemma 40.** *Assume that  $S$  is a good set. We have*

$$\sup_{p_1 \in \mathcal{P}_{n,d,2k}^1} \left| \mathbb{E}[p_1^2(S)] - 1 \right| \leq 2 \cdot \beta(2\delta_\gamma, d, \sqrt{2}) + \alpha_1 \gamma_1^2.$$

In particular, as we set  $\alpha_1 \leq \frac{\epsilon}{\gamma_1^2}$ , we have

$$\sup_{p_1 \in \mathcal{P}_{n,d,2k}^1} \left| \mathbb{E}[p_1^2(S)] - 1 \right| \leq 2 \cdot \beta(2\delta_\gamma, d, \sqrt{2}) + \epsilon.$$

*Proof.* By Fact 20,

$$\begin{aligned} \left| \mathbb{E}[p_1^2(D)] - \mathbb{E}[p_1^2(D|_{X_\gamma})] \right| &= \left| \int_0^\infty 2t \left[ \Pr(|p_1(D)| > t) - \Pr(|p_1(D|_{X_\gamma})| > t) \right] dt \right| \\ &\leq \int_0^\infty 2t \min\{2\delta_\gamma, \Pr(|p_1(D)| \geq t)\} dt \\ &\leq 2\beta(2\delta_\gamma, d, \sqrt{2}), \end{aligned}$$

where the second step follows from Lemma 44 and the last step follows from Lemma 24.

On the other hand, by Part 3 of Definition 12, we have

$$\begin{aligned} \left| \mathbb{E}[p_1^2(S|_{X_\gamma})] - \mathbb{E}[p_1^2(D|_{X_\gamma})] \right| &= \left| \int_0^\infty 2t \left[ \Pr(|p_1(S|_{X_\gamma})| > t) - \Pr(|p_1(D|_{X_\gamma})| > t) \right] dt \right| \\ &\leq \int_0^{\gamma_1} 2t \alpha_1 dt \\ &= \alpha_1 \gamma_1^2, \end{aligned}$$

where the inequality follows since  $|p_1(x)| \leq \gamma_1$  for all  $x \in X_\gamma$  (Lemma 29). By triangle inequality, the fact that  $\mathbb{E}[p_1^2(D)] = 1$  (Lemma 29), and  $S|_{X_\gamma} = S$  ( $S$  is a good set), we complete the proof.  $\square$

**Lemma 41.** *Assume that  $S$  is a good set and  $\Delta(S, S') \leq 2\eta$ . Let  $L = S \setminus S'$ . We have*

$$\sup_{p_1 \in \mathcal{P}_{n,d,2k}^1} (|L| / |S|) \cdot \mathbb{E}[p_1^2(L)] \leq 12 \cdot \beta(\eta, d, \sqrt{2}) + 4\eta + \alpha_1 \gamma_1^2.$$

In particular, as we set  $\alpha_1 \leq \frac{\epsilon}{\gamma_1^2}$ , we have

$$\sup_{p_1 \in \mathcal{P}_{n,d,2k}^1} (|L| / |S|) \cdot \mathbb{E}[p_1^2(L)] \leq 12 \cdot \beta(\eta, d, \sqrt{2}) + 4\eta + \epsilon.$$

*Proof.* We will use Lemma 43 to show the result. Since  $S$  is a good set, we have  $S|_{X_\gamma} = S$ . We have  $|L| \leq 2\eta|S| = |S|$  since  $\Delta(S, S') \leq 2\eta$ . By Lemma 29 and Part 2

in that lemma, we set  $S_0 = S$ ,  $S_1 = L$ ,  $\omega_1 = 2$ ,  $\epsilon_0 = \eta$ . By Lemma 29, we set  $Q_d(t) = e^{-(t/\sqrt{2})^{2/d}/c_0}$ , and  $t_0 = \sqrt{2}$ . Lemma 29 tells that we can set  $\omega_2 = 1$ . In this way, by Corollary 25, we set  $\beta_0 = \beta(\eta, d, \sqrt{2})$ . By Lemma 29, we can set  $\gamma_0 = \gamma_1$ . Therefore, we obtain the desired bound.  $\square$

**Lemma 42.** *Assume that  $S$  is a good set,  $\Delta(S, S') \leq 2\eta$ , and  $\|(\Sigma - I)_U\|_F \leq \kappa$ . We have*

$$\sup_{p_1 \in \mathcal{P}_{n,d,2k}^1} |E| \cdot \mathbb{E}[p_1^2(E)] \leq |S'| \cdot \frac{6}{c} (\kappa + \beta(\eta, d, \sqrt{2}) + \beta(2\delta_\gamma, d, \sqrt{2}) + \eta + \epsilon).$$

*Proof.* Recall  $S' = S \cup E \setminus L$ . By algebraic calculation, we have

$$\begin{aligned} |E| \cdot \mathbb{E}[p_1^2(E)] &= |S'| \cdot \mathbb{E}[p_1^2(S')] + |L| \cdot \mathbb{E}[p_1^2(L)] - |S| \cdot \mathbb{E}[p_1^2(S)] \\ &\leq |S'| \cdot (\kappa + 1) + |S| \cdot (12\beta(\eta, d, \sqrt{2}) + 4\eta + \epsilon) - |S| \cdot (1 - 2 \cdot \beta(2\delta_\gamma, d, \sqrt{2}) - \epsilon) \\ &\leq |S'| \cdot \kappa + 12|S| (\beta(\eta, d, \sqrt{2}) + \beta(2\delta_\gamma, d, \sqrt{2}) + \eta + \epsilon) \end{aligned}$$

where we applied Lemma 39, Lemma 40 and Lemma 41 in the first inequality and the fact  $|S'| \leq |S|$  in the last step. Since  $\Delta(S, S') \leq 2\eta$ , in view of (A.2), we have  $|S| \leq \frac{1}{1-2\eta} |S'| \leq \frac{1}{2c} |S'|$ . Rearranging gives the desired result.  $\square$

The following lemma is similar to Lemma 37, but we upper bound the expectation of the square of a polynomial.

**Lemma 43.** *Let  $c_0 > 0$  be an absolute constant. Let  $S_0$  be a set of instances in  $\mathbb{R}^n$  and  $S_1 \subset S_0$ , with  $|S_1| \leq \omega_1 \tau |S_0|$  for  $\omega_1, \tau > 0$ . Let  $p$  be such that  $\Pr(|p(S_0)| > t) \leq \omega_2 \cdot Q_d(t) + \alpha_0$  for all  $t \geq t_0$ , where  $\omega_2, Q_d(t), \alpha_0 > 0$ . Assume  $\max_{x \in S_0} |p(x)| \leq \gamma_0$ . Further assume that  $\int_0^\infty t \min\{\tau, Q_d(t)\} dt \leq \beta_0$ . Then*

$$(|S_1| / |S_0|) \cdot \mathbb{E}[p^2(S_1)] \leq \omega_1 t_0^2 \cdot \tau + 2(\omega_1 + 1)(\omega_2 + 1)\beta_0 + \alpha_0 \gamma_0^2.$$

*Proof.* Since  $S_1 \subset S_0$ , we have  $|S_1| \cdot \Pr(|p(S_1)| > t) \leq |S_0| \cdot \Pr(|p(S_0)| > t)$ . Thus,

$$\Pr(|p(S_1)| > t) \leq \min \left\{ 1, \frac{|S_0|}{|S_1|} \cdot \Pr(|p(S_0)| > t) \right\}. \quad (\text{E.10})$$

By Fact 20, we have

$$\begin{aligned} (|S_1| / |S_0|) \cdot \mathbb{E}[p^2(S_1)] &\leq \int_0^\infty 2t(|S_1| / |S_0|) \Pr(|p(S_1)| > t) dt \\ &\stackrel{\zeta_1}{=} \int_0^{\gamma_0} 2t(|S_1| / |S_0|) \Pr(|p(S_1)| > t) dt \\ &\stackrel{\zeta_2}{\leq} \int_0^{\gamma_0} 2t \min \left\{ |S_1| / |S_0|, \Pr(|p(S_0)| > t) \right\} dt \\ &\stackrel{\zeta_3}{\leq} \int_0^{\gamma_0} 2t \min \left\{ \omega_1 \tau, \Pr(|p(S_0)| > t) \right\} dt \\ &\stackrel{\zeta_4}{\leq} \int_0^{t_0} 2t \min \{\omega_1 \tau, 1\} dt + \int_{t_0}^{\gamma_0} 2t \min \{\omega_1 \tau, \omega_2 \cdot Q_d(t) + \alpha_0\} dt \\ &\stackrel{\zeta_5}{\leq} \int_0^{t_0} 2\omega_1 \tau t dt + \int_{t_0}^{\gamma_0} 2t \min \{\omega_1 \tau, \omega_2 \cdot Q_d(t)\} dt + \int_{t_0}^{\gamma_0} 2t \alpha_0 dt \\ &\stackrel{\zeta_6}{\leq} \omega_1 t_0^2 \cdot \tau + 2(\omega_1 + 1)(\omega_2 + 1) \int_{t_0}^{\gamma_0} t \min \{\tau, Q_d(t)\} dt + \alpha_0 (\gamma_0^2 - t_0^2) \\ &\stackrel{\zeta_7}{\leq} \omega_1 t_0^2 \cdot \tau + 2(\omega_1 + 1)(\omega_2 + 1)\beta_0 + \alpha_0 \gamma_0^2. \end{aligned}$$

In the above,  $\zeta_1$  follows from the condition that  $p(x) \leq \gamma_0$  for all  $x \in S_1$ ,  $\zeta_2$  follows from (E.10),  $\zeta_3$  uses the condition  $|S_1| \leq \omega_1 \tau |S_0|$ ,  $\zeta_4$  uses the condition of the tail bound of  $p(S_0)$  when  $t \geq t_0$ ,  $\zeta_5$  applies elementary facts that  $\min\{\omega_1 \tau, 1\} \leq \omega_1 \tau$  and  $\min\{a, b+c\} \leq \min\{a, b\} + c$  for any  $c > 0$ ,  $\zeta_6$  uses the fact that both  $\frac{\omega_1}{(\omega_1+1)(\omega_2+1)}$  and  $\frac{\omega_1}{(\omega_1+1)(\omega_2+1)}$  are less than 1 for positive  $\omega_1$  and  $\omega_2$ , and  $\zeta_7$  applies the condition on the integral and uses the fact that both  $\tau$  and  $Q_d(t)$  are positive.  $\square$

**Lemma 44.** Suppose  $\delta_\gamma \leq 1/2$ . The following holds for any function  $p$ :

$$\left| \Pr(|p(D|_{X_\gamma})| \geq t) - \Pr(|p(D)| \geq t) \right| \leq \min \{2\delta_\gamma, \Pr(|p(D)| \geq t)\}.$$

*Proof.* Lemma 11 tells that  $\Pr_{x \sim D}(x \notin X_\gamma) \leq \delta_\gamma$ . Observe that

$$\Pr(|p(D|_{X_\gamma})| \geq t) \leq \frac{\Pr(|p(D)| \geq t)}{\Pr_{x \sim D}(x \in X_\gamma)} \leq \frac{1}{1 - \delta_\gamma} \Pr(|p(D)| \geq t) \leq 2 \Pr(|p(D)| \geq t),$$

implying

$$\left| \Pr(|p(D|_{X_\gamma})| \geq t) - \Pr(|p(D)| \geq t) \right| \leq \Pr(|p(D)| \geq t). \quad (\text{E.11})$$

On the other hand, for any event  $A$ ,

$$\Pr_D(A) = \Pr_D(A \mid x \in X_\gamma) \cdot \Pr_D(x \in X_\gamma) + \Pr_D(A \mid x \notin X_\gamma) \cdot \Pr_D(x \notin X_\gamma).$$

This implies

$$\begin{aligned} & |\Pr_D(A) - \Pr_D(A \mid x \in X_\gamma)| \\ &= |\Pr_D(A \mid x \in X_\gamma) \cdot (\Pr_D(x \in X_\gamma) - 1) + \Pr_D(A \mid x \notin X_\gamma) \cdot \Pr_D(x \notin X_\gamma)| \\ &= |- \Pr_D(A \mid x \in X_\gamma) \cdot \Pr_D(x \notin X_\gamma) + \Pr_D(A \mid x \notin X_\gamma) \cdot \Pr_D(x \notin X_\gamma)| \\ &\leq 2 \Pr_D(x \notin X_\gamma) \\ &\leq 2\delta_\gamma. \end{aligned}$$

This completes the proof.  $\square$

**Lemma 45.** For any vector  $w$  and any  $k$ -sparse vector  $u$ , if  $\sup_{v: \|v\|_2=1, \|v\|_0 \leq 2k} |\langle v, w - u \rangle| \leq \epsilon$ , then  $\|\mathbf{H}_k(w) - u\|_2 \leq 3\epsilon$ .

*Proof.* Let  $\Lambda_0$  be the support set of  $\mathbf{H}_k(w)$ ,  $\Lambda_1 = \text{supp}(u) \setminus \Lambda_0$ ,  $\Lambda_2 = \text{supp}(u) \cap \Lambda_0$ ,  $\Lambda_3 = \Lambda_0 \setminus \text{supp}(u)$ . Therefore, we can decompose

$$\|\mathbf{H}_k(w) - u\|_2^2 = \|u_{\Lambda_1}\|_2^2 + \|(w - u)_{\Lambda_2}\|_2^2 + \|w_{\Lambda_3}\|_2^2. \quad (\text{E.12})$$

Note that by choosing  $v = (w - u)_{\Lambda_2 \cup \Lambda_3} / \|(w - u)_{\Lambda_2 \cup \Lambda_3}\|_2$ , we get

$$\|(w - u)_{\Lambda_2}\|_2^2 + \|w_{\Lambda_3}\|_2^2 \leq \epsilon^2. \quad (\text{E.13})$$

On the other side, observe that

$$\|u_{\Lambda_1}\|_2 \leq \|(u - w)_{\Lambda_1}\|_2 + \|w_{\Lambda_1}\|_2 \quad (\text{E.14})$$

by triangle inequality. Since  $\Lambda_3$  is a subset of  $\Lambda_0$ , the index set of the  $k$  largest elements of  $w$ , while  $\Lambda_1 \cap \Lambda_0 = \emptyset$ , we know that elements of  $w$  in  $\Lambda_1$  are less than those in  $\Lambda_3$ . This combined with the fact that  $|\Lambda_1| = |\Lambda_3|$  implies that

$$\|w_{\Lambda_1}\|_2 \leq \|w_{\Lambda_3}\|_2 \leq \epsilon. \quad (\text{E.15})$$

where the second step follows from (E.13). In order to upper bound  $\|(u - w)_{\Lambda_1}\|_2$ , we can pick  $v = (u - w)_{\Lambda_1} / \|(u - w)_{\Lambda_1}\|_2$  and get

$$\|(u - w)_{\Lambda_1}\|_2 \leq \epsilon. \quad (\text{E.16})$$

Plugging (E.15) and (E.16) into (E.14) gives

$$\|u_{\Lambda_1}\|_2 \leq 2\epsilon.$$

This in conjunction with (E.13) and (E.12) gives  $\|\mathbf{H}_k(w) - u\|_2 \leq \sqrt{5}\epsilon \leq 3\epsilon$ .  $\square$

## F. Proof of Theorem 19

*Proof.* The sample complexity and the first equation are an immediate result from Theorem 18 and Lemma 46. The second equation follows from algebraic calculation.  $\square$

**Lemma 46** (Diakonikolas et al. (2018a)). *Suppose  $D$  is  $\mathcal{N}(0, I_{n \times n})$ . There is an algorithm that takes as input a vector  $u \in \mathbb{R}^{(n+1)^d}$  with  $\|u - \chi_{f^*}\|_2 \leq \epsilon$ , runs in time  $O(n^d/\epsilon^2)$  and outputs a degree- $d$  PTF  $\hat{f}$  such that*

$$\Pr_{x \sim D} (\hat{f}(x) \neq f^*(x)) \leq c_1 \cdot d \cdot \epsilon^{\frac{1}{d+1}},$$

for some absolute constant  $c_1 > 0$ .

*Proof.* The result is a combination of Theorem 10 of De et al. (2014), Lemma 3.4 and Lemma 3.5 of Diakonikolas et al. (2018a). The only difference is that when applying Theorem 10 of De et al. (2014) to our setup, we can compute Chow vectors of a given function exactly since  $D$  is known to be Gaussian; thus no additional samples are needed and the running time is slightly better than their original analysis. See Lemma 47 for clarification.  $\square$

We reproduce the proof of Theorem 10 of De et al. (2014) but tailored to our case that  $D$  is Gaussian, and thus there is no need to acquire additional samples.

**Lemma 47.** *Let  $f$  be a degree- $d$  PTF on  $\mathbb{R}^n$ . There is an algorithm that takes as input a vector  $v$  with  $\|v - \chi_f\|_2 \leq \epsilon$ , and outputs a polynomial bounded function  $g : \mathbb{R}^n \rightarrow [-1, 1]$  such that  $\|\chi_g - \chi_f\|_2 \leq 4\epsilon$ . In addition, the algorithm runs in  $O(n^d/\epsilon^2)$  time.*

*Proof.* We will iteratively construct a sequence of functions  $\{g_t\} \subset \mathcal{P}_{n,d}$ . Let  $g'_0 = 0$  and  $g_0 = P_1(g'_0)$ , where  $P_1(a) = \text{sign}(a)$  if  $|a| \geq 1$  and  $P_1(a) = a$  otherwise. Given  $g_t$ , we compute  $\chi_{g_t}$ . Let

$$\rho := \|v - \chi_{g_t}\|_2. \quad (\text{F.1})$$

**Case 1.**  $\rho \leq 3\epsilon$ . Then

$$\|\chi_{g_t} - \chi_f\|_2 \leq \|\chi_{g_t} - v\|_2 + \|v - \chi_f\|_2 = \|\chi_{g_t} - v\|_2 + \|v - \chi_f\|_2 \leq 4\epsilon.$$

Thus we output  $g_t$ .

**Case 2.**  $\rho > 3\epsilon$ . Define

$$h'_t(x) = (v - \chi_{g_t}) \cdot m(x), \quad g'_{t+1} = g'_t + \frac{1}{2}h'_t, \quad g_{t+1} = P_1(g'_{t+1}). \quad (\text{F.2})$$

Consider the following potential function:

$$E(t) = \mathbb{E}[(f - g_t)(f - 2g'_t + g_t)]. \quad (\text{F.3})$$

The proof of Theorem 10 of De et al. (2014) implies the following two claims:

**Claim 48.**  $\mathbb{E}[(f - g_t)h'_t] \geq \rho(\rho - \epsilon)$ .

**Claim 49.** Given any  $g'_t$  and  $h'_t$ , let  $g_t = P_1(g'_t)$ ,  $g'_{t+1} = g'_t + \frac{1}{2}h'_t$ ,  $g_{t+1} = P_1(g'_{t+1})$ . Then  $\mathbb{E}[(g_{t+1} - g_t)(2g'_{t+1} - g_t - g_{t+1})] \leq \frac{1}{2}\mathbb{E}[(h'_t)^2]$ .

Observe that by our definition of  $h'_t$ , we have

$$\mathbb{E}[(h'_t)^2] = \|v - \chi_{g_t}\|_2^2 = \rho^2. \quad (\text{F.4})$$

Therefore, the progress of  $E(t)$  can be bounded as follows:

$$\begin{aligned}
 E(t+1) - E(t) &= -\mathbb{E}[(f - g_t)h'_t] + \mathbb{E}[(g_{t+1} - g_t)(2g'_{t+1} - g_t - g_{t+1})] \\
 &\leq -\rho(\rho - \epsilon) + \frac{1}{2}\rho^2 \\
 &\leq -\epsilon^2.
 \end{aligned}$$

In addition, we have  $E(t) \geq 0$  and  $E(0) = 1$ . These together imply that the algorithm terminates in at most  $\frac{1}{\epsilon^2}$  iterations. It is easy to see that the computational cost in each iteration is dominated by the construction of  $h'_t(\cdot)$ , which is  $O(n^d)$ . Thus, the overall running time is  $O(n^d/\epsilon^2)$ .  $\square$