

# AN EFFICIENT TESTER-LEARNER FOR HALFSPACES

Aravind Gollakota\* Adam R. Klivans† Konstantinos Stavropoulos‡ Arsen Vasilyan§  
 Apple UT Austin UT Austin MIT

## ABSTRACT

We give the first efficient algorithm for learning halfspaces in the testable learning model recently defined by Rubinfeld and Vasilyan (RV23). In this model, a learner certifies that the accuracy of its output hypothesis is near optimal whenever the training set passes an associated test, and training sets drawn from some target distribution must pass the test. This model is more challenging than distribution-specific agnostic or Massart noise models where the learner is allowed to fail arbitrarily if the distributional assumption does not hold. We consider the setting where the target distribution is the standard Gaussian in  $d$  dimensions and the label noise is either Massart or adversarial (agnostic). For Massart noise, our tester-learner runs in polynomial time and outputs a hypothesis with (information-theoretically optimal) error  $\text{opt} + \epsilon$  (and extends to any fixed strongly log-concave target distribution). For adversarial noise, our tester-learner obtains error  $O(\text{opt}) + \epsilon$  in polynomial time. Prior work on testable learning ignores the labels in the training set and checks that the empirical moments of the covariates are close to the moments of the base distribution. Here we develop new tests of independent interest that make critical use of the labels and combine them with the moment-matching approach of (GKK23). This enables us to implement a testable variant of the algorithm of (DKTZ20a; DKTZ20b) for learning noisy halfspaces using nonconvex SGD.

## 1 INTRODUCTION

Learning halfspaces in the presence of noise is one of the most basic and well-studied problems in computational learning theory. A large body of work has obtained results for this problem under a variety of different noise models and distributional assumptions (see e.g. (BH21) for a survey). A major issue with common distributional assumptions such as Gaussianity, however, is that they can be hard or impossible to verify in the absence of any prior information.

The recently defined model of testable learning (RV23) addresses this issue by replacing such assumptions with efficiently testable ones. In this model, the learner is required to work with an arbitrary input distribution  $D_{\mathcal{X}\mathcal{Y}}$  and verify any assumptions it needs to succeed. It may choose to reject a given training set, but if it accepts, it is required to output a hypothesis with error close to  $\text{opt}(\mathcal{C}, D_{\mathcal{X}\mathcal{Y}})$ , the optimal error achievable over  $D_{\mathcal{X}\mathcal{Y}}$  by any function in a concept class  $\mathcal{C}$ . Further, whenever the training set is drawn from a distribution  $D_{\mathcal{X}\mathcal{Y}}$  whose marginal is truly a well-behaved target distribution  $D^*$  (such as the standard Gaussian), the algorithm is required to accept with high probability. Such an algorithm, or tester-learner, is then said to testably learn  $\mathcal{C}$  with respect to target marginal  $D^*$ . (See Definition 2.1.) Note that unlike ordinary distribution-specific agnostic learners, a tester-learner must take some nontrivial action *regardless* of the input distribution.

The work of (RV23; GKK23) established foundational algorithmic and statistical results for this model and showed that testable learning is in general provably harder than ordinary distribution-specific agnostic learning. As one of their main algorithmic results, they showed tester-learners for the class of halfspaces over  $\mathbb{R}^d$  that succeed whenever the target marginal is Gaussian (or one of a more general class of distributions), achieving error  $\text{opt} + \epsilon$  in time and sample complexity

\*aravindg@cs.utexas.edu

†klivans@cs.utexas.edu

‡kstavrop@cs.utexas.edu

§vasilyan@mit.edu

$d^{\Theta(1/\epsilon^2)}$ . This matches the running time of ordinary distribution-specific agnostic learning of halfspaces over the Gaussian using the standard approach of (KKMS08). Their testers are simple and label-oblivious, and are based on checking whether the low-degree empirical moments of the unknown marginal match those of the target  $D^*$ .

These works essentially resolve the question of designing tester-learners achieving error  $\text{opt} + \epsilon$  for halfspaces, matching known hardness results for (ordinary) agnostic learning (GGK20; DKZ20; DKPZ21). Their running time, however, necessarily scales exponentially in  $1/\epsilon$ .

A long line of research has sought to obtain more efficient algorithms at the cost of relaxing the optimality guarantee (ABL17; DKS18; DKTZ20a; DKTZ20b). These works give polynomial-time algorithms achieving bounds of the form  $\text{opt} + \epsilon$  and  $O(\text{opt}) + \epsilon$  for the Massart and agnostic setting respectively under structured distributions (see Section 1.1 for more discussion). The main question we consider here is whether such guarantees can be obtained in the testable learning framework.

**Our contributions.** In this work we design the first tester-learners for halfspaces that run in fully polynomial time in all parameters. We match the optimality guarantees of fully polynomial-time learning algorithms under Gaussian marginals for the Massart noise model (where the labels arise from a halfspace but are flipped by an adversary with probability at most  $\eta$ ) as well as for the agnostic model (where the labels can be completely arbitrary). In fact, for the Massart setting our guarantee holds with respect to any chosen target marginal  $D^*$  that is isotropic and strongly log-concave, and the same is true of the agnostic setting albeit with a slightly weaker guarantee.

**Theorem 1.1** (Formally stated as Theorem 4.1). *Let  $\mathcal{C}$  be the class of origin-centered halfspaces over  $\mathbb{R}^d$ , and let  $D^*$  be any isotropic strongly log-concave distribution. In the setting where the labels are corrupted with Massart noise at rate at most  $\eta < \frac{1}{2}$ ,  $\mathcal{C}$  can be testably learned w.r.t.  $D^*$  up to error  $\text{opt} + \epsilon$  using  $\text{poly}(d, \frac{1}{\epsilon}, \frac{1}{1-2\eta})$  time and sample complexity.*

**Theorem 1.2** (Formally stated as Theorem 5.1). *Let  $\mathcal{C}$  be as above. In the adversarial noise or agnostic setting where the labels are completely arbitrary,  $\mathcal{C}$  can be testably learned w.r.t.  $\mathcal{N}(0, I_d)$  up to error  $O(\text{opt}) + \epsilon$  using  $\text{poly}(d, \frac{1}{\epsilon})$  time and sample complexity.*

**Our techniques.** The tester-learners we develop are significantly more involved than prior work on testable learning. We build on the nonconvex optimization approach to learning noisy halfspaces due to (DKTZ20a; DKTZ20b) as well as the structural results on fooling functions of halfspaces using moment matching due to (GKK23). Unlike the label-oblivious, global moment tests of (RV23; GKK23), our tests make crucial use of the labels and check *local* properties of the distribution in regions described by certain candidate vectors. These candidates are approximate stationary points of a natural nonconvex surrogate of the 0-1 loss, obtained by running gradient descent. When the distribution is known to be well-behaved, (DKTZ20a; DKTZ20b) showed that any such stationary point is in fact a good solution (for technical reasons we must use a slightly different surrogate loss). Their proof relies crucially on structural geometric properties that hold for these well-behaved distributions, an important one being that the probability mass of any region close to the origin is proportional to its geometric measure.

In the testable learning setting, we must efficiently check this property for candidate solutions. Since these regions may be described as intersections of halfspaces, we may hope to apply the moment-matching framework of (GKK23). Naïvely, however, they only allow us to check in polynomial time that the probability masses of such regions are within an additive constant of what they should be under the target marginal. But we can view these regions as sub-regions of a known band described by our candidate vector. By running moment tests on the distribution *conditioned* on this band and exploiting the full strength of the moment-matching framework, we are able to effectively convert our weak additive approximations to good multiplicative ones. This allows us to argue that our stationary points are indeed good solutions.

**Independent and Subsequent Works.** In this paper we provide the first efficient tester-learners for halfspaces when the noise is either adversarial or Massart. In independent and concurrent work by (DKK<sup>+</sup>23), an efficient tester-learner for homogeneous halfspaces achieving error  $O(\text{opt}) + \epsilon$  for Gaussian target marginals is also provided, but they do not provide any results for arbitrary strongly log-concave target marginals (see Theorem 5.1) or a guarantee of  $\text{opt} + \epsilon$  for Massart noise. In subsequent work by (GKSV23), our techniques were used to provide tester-learners that are not tailored to a single target distribution, but are guaranteed to accept any member of a large family of distributions. Although their main results are more general, their approach crucially extends our

approach here. Moreover, on the technical side, the proof we give here shows how to make use of the moment-matching approach of (GKK23) to provide fully polynomial-time efficient tester-learners, which might be of independent interest.

## 1.1 RELATED WORK

We provide a partial summary of some of the most relevant prior and related work on efficient algorithms for learning halfspaces in the presence of adversarial label or Massart noise, and refer the reader to (BH21) for a survey.

In the distribution-specific agnostic setting where the marginal is assumed to be isotropic and log-concave, (KLS09) showed an algorithm achieving  $\epsilon \sqrt{\log(1/\epsilon)}$  for the class of origin-centered halfspaces. (ABL17) later obtained  $\epsilon \sqrt{\log(1/\epsilon)}$  using an approach that introduced the principle of iterative localization, where the learner focuses attention on a band around a candidate halfspace in order to produce an improved candidate. (Dan15) used this principle to obtain a PTAS for agnostically learning halfspaces under the uniform distribution on the sphere, and (BZ17) extended it to more general s-concave distributions. Further works in this line include (YZ17; Zha18; ZSA20; ZL21). (DKTZ20b) introduced the simplest approach yet, based entirely on nonconvex SGD, and showed that it achieves  $\epsilon \sqrt{\log(1/\epsilon)}$  for origin-centered halfspaces over a wide class of structured distributions. Other related works include (DKS18; DKTZ22).

In the Massart noise setting with noise rate bounded by  $\eta$ , work of (DGT19) gave the first efficient distribution-free algorithm achieving error  $\epsilon \sqrt{\log(1/\epsilon)}$ ; further improvements and followups include (DKT21; DTK22). However, the optimal error  $\epsilon \sqrt{\log(1/\epsilon)}$  achievable by a halfspace may be much smaller than  $\epsilon \sqrt{\log(1/\epsilon)}$ , and it has been shown that there are distributions where achieving error competitive with  $\epsilon \sqrt{\log(1/\epsilon)}$  as opposed to is computationally hard (DK22; DKMR22). As a result, the distribution-specific setting remains well-motivated for Massart noise. Early distribution-specific algorithms were given by (ABHU15; ABHZ16), but a key breakthrough was the nonconvex SGD approach introduced by (DKTZ20a), which achieved error  $\epsilon \sqrt{\log(1/\epsilon)}$  for origin-centered halfspaces efficiently over a wide range of distributions. This was later generalized by (DK22).

## 1.2 TECHNICAL OVERVIEW

Our starting point is the nonconvex optimization approach to learning noisy halfspaces due to (DKTZ20a; DKTZ20b). The algorithms in these works consist of running SGD on a natural nonconvex surrogate for the 0-1 loss, namely a smooth version of the ramp loss. The key structural property shown is that if the marginal distribution is structured (e.g. log-concave) and the slope of the ramp is picked appropriately, then any  $w$  that has large angle with an optimal  $w^*$  cannot be an approximate stationary point of the surrogate loss, i.e.  $\| \nabla L(w) \|$  must be large. This is proven by carefully analyzing the contributions to the gradient norm from certain critical regions of  $\text{span}(w; w^*)$ , and crucially using the distributional assumption that the probability masses of these regions are proportional to their geometric measures. (See Fig. 3.) In the testable learning setting, the main challenge we face in adapting this approach is checking such a property for the unknown distribution we have access to.

A preliminary observation is that the critical regions  $\text{span}(w; w^*)$  that we need to analyze are rectangles, and are hence functions of a small number of halfspaces. Encouragingly, one of the key structural results of the prior work of (GKK23) pertains to “fooling” such functions. Concretely, they show that whenever the true marginal  $D$  matches moments of degree at most  $t = 2$  with a target  $D^*$  that satisfies suitable concentration and anticoncentration properties, then  $\mathbb{E}_D[f] \approx \mathbb{E}_{D^*}[f]$  for any  $f$  that is a function of a small number of halfspaces. If we could run such a test and ensure that the probabilities of the critical regions over our empirical marginal are also related to their areas, then we would have a similar stationary point property. However, the difficulty is that since we wish to run in fully polynomial time, we can only hope to fool such functions up to a constant. Unfortunately, this is not sufficient to analyze the probability masses of the critical regions we care about as they may be very small.

The chief insight that lets us get around this issue is that each critical region is in fact of a very specific form, namely a rectangle that is axis-aligned with  $R = \{x : |x_i - x_j| \leq \delta_i, |x_i - x_j| \leq \delta_j\}$  for some values  $\delta_i, \delta_j$  and some  $v$  orthogonal to  $w$ . Moreover, we know  $w$ , meaning

we can efficiently estimate the probability  $\mathbb{P}_{D_X}[\langle w; x_i \rangle \geq \tau]$  up to constant multiplicative factors without needing moment tests. Denoting the band  $\langle w; x_i \rangle \geq \tau$  by  $T$  and writing  $\mathbb{P}_{D_X}[R] = \mathbb{P}_{D_X}[\langle w; x_i \rangle \geq \tau] \times \mathbb{P}_{D_X}[T]$ , it turns out that we should expect  $\mathbb{P}_{D_X}[\langle w; x_i \rangle \geq \tau] = (1 - \epsilon)$ , as this is what would occur under the structured target distribution  $D$ . (Such a “localization” property is also at the heart of the algorithms for approximately learning halfspaces of, e.g., (ABL17; Dan15).) To check this, it suffices to run tests that ensure that  $\mathbb{P}_{D_X}[\langle w; x_i \rangle \geq \tau]$  is within an additive constant of this probability under  $D$ .

We can now describe the core of our algorithm (omitting some details such as the selection of the slope of the ramp). First, we run SGD on the surrogate loss and arrive at an approximate stationary point and candidate vector (technically a list of such candidates). Then, we define the band  $T$  based on  $w$ , and run tests on the empirical distribution conditioned on  $T$ . Specifically, we check that the low-degree empirical moments conditioned on  $T$  match those of  $D$  conditioned on  $T$ , and then apply the structural result of (GKK23) to ensure conditional probabilities of the form  $\mathbb{P}_{D_X}[\langle w; x_i \rangle \geq \tau | x_i \in T]$  match  $\mathbb{P}_D[\langle w; x_i \rangle \geq \tau | x_i \in T]$  up to a suitable additive constant. This suffices to ensure that even over our empirical marginal, the particular stationary point we have is indeed close in angular distance to an optimal

Angular distance that remains, often taken for granted under structured distributions, is that closeness in angular distance  $\angle(w; w^*)$  does not immediately translate to closeness in terms of agreement,  $\mathbb{P}[\text{sign}(w; x_i) \neq \text{sign}(w^*; x_i)]$ , over our unknown marginal. Nevertheless, we show that when the target distribution is Gaussian, we can run polynomial-time tests that ensure that an angle of  $\angle(w; w^*) = \epsilon$  translates to disagreement of at most  $\epsilon$ . When the target distribution is a general strongly log-concave distribution, we show a slightly weaker relationship: for any  $\epsilon$ , we can run tests requiring time  $d^{\epsilon^{-k}}$  that ensure that an angle of  $\epsilon$  translates to disagreement of at most  $O(\epsilon^{-k-1})$ . In the Massart noise setting, we can make  $\angle(w; w^*)$  arbitrarily small, and so obtain our  $\epsilon$ -guarantee for any target strongly log-concave distribution in polynomial time. In the adversarial noise setting, we face a more delicate tradeoff and can only make  $\angle(w; w^*)$  as small as  $\epsilon^{\text{opt}}$ . When the target distribution is Gaussian, this is enough to obtain  $\epsilon^{\text{opt}}$  in polynomial time. When the target distribution is a general strongly log-concave distribution, we instead obtain  $\epsilon^{\text{opt}}$  in quasipolynomial time.

## 2 PRELIMINARIES

**Notation and setup** Throughout, the domain will be  $X = \mathbb{R}^d$ , and labels will lie in  $Y = \{-1, 1\}$ . The unknown joint distribution over  $X \times Y$  that we have access to will be denoted  $D_{XY}$ , and its marginal on  $X$  will be denoted by  $D_X$ . The target marginal on  $X$  will be denoted by  $D$ . We use the following convention for monomials: for a multi-index  $\alpha = (\alpha_1, \dots, \alpha_d) \in \mathbb{Z}_{\geq 0}^d$ ,  $x^\alpha$  denotes  $\prod_i x_i^{\alpha_i}$ , and  $|\alpha| = \sum_i \alpha_i$  denotes its total degree. We use  $\mathcal{C}$  to denote a concept class mapping  $\mathbb{R}^d$  to  $\{-1, 1\}$ , which throughout this paper will be the class of halfspaces or functions of halfspaces over  $\mathbb{R}^d$ . We use  $\text{opt}(\mathcal{C}; D_{XY})$  to denote the optimal error  $\inf_{f \in \mathcal{C}} \mathbb{E}_{(x,y) \sim D_{XY}} [f(x) \neq y]$ , or just  $\text{opt}$  when  $\mathcal{C}$  and  $D_{XY}$  are clear from context. We recall the definitions of the noise models we consider. In the Massart noise model, the labels satisfy  $\mathbb{P}_{D_{XY}}[y \neq \text{sign}(\langle w; x_i \rangle)] = \epsilon(x)$ , where  $\epsilon(x) < \frac{1}{2}$  for all  $x$ . In the adversarial label noise or agnostic model, the labels may be completely arbitrary. In both cases, the learner’s goal is to produce a hypothesis with error competitive with  $\text{opt}$ . We now formally define testable learning. The following definition is an equivalent reframing of the original definition (RV23, Def 4), folding the (label-aware) tester and learner into a single tester-learner.

**Definition 2.1 (Testable learning, (RV23))** Let  $\mathcal{C}$  be a concept class mapping  $\mathbb{R}^d$  to  $\{-1, 1\}$ . Let  $D$  be a certain target marginal on  $\mathbb{R}^d$ . Let  $\epsilon, \delta > 0$  be parameters, and let  $\gamma: [0, 1] \rightarrow [0, 1]$  be some function. We say  $\mathcal{C}$  can be testably learned w.r.t.  $D$  up to error  $\gamma(\text{opt}) + \delta$  with failure probability  $\delta$  if there exists a tester-learner  $\mathcal{A}$  meeting the following specification. For any distribution  $D_{XY}$  on  $\mathbb{R}^d \times \{-1, 1\}$ ,  $\mathcal{A}$  takes in a large sample drawn from  $D_{XY}$ , and either rejects or accepts and produces a hypothesis  $h: \mathbb{R}^d \rightarrow \{-1, 1\}$ . Further, the following conditions must be met:

- (Soundness.) Whenever  $\mathcal{A}$  accepts and produces a hypothesis  $h$ , with probability at least  $1 - \delta$  (over the randomness of  $\mathcal{A}$  and  $D_{XY}$ ),  $h$  must satisfy  $\mathbb{E}_{(x,y) \sim D_{XY}} [h(x) \neq y] \leq \gamma(\text{opt}(\mathcal{C}; D_{XY})) + \delta$ .

- (b) (Completeness.) Whenever  $\mathcal{D}_{X,Y}$  truly has marginal  $\mathcal{D}$ ,  $A$  must accept with probability at least  $1 - \epsilon$  (over the randomness of  $\mathcal{D}$  and  $A$ ).

### 3 TESTING PROPERTIES OF STRONGLY LOG-CONCAVE DISTRIBUTIONS

In this section we define the testers that we will need for our algorithm. All the proofs from this section can be found in Appendix B. We begin with a structural lemma that strengthens the key structural result of (GKK23), stated here as Proposition A.3. It states that even when we restrict an isotropic strongly log-concave  $\mathcal{D}$  to a band around the origin, moment matching suffices to fool functions of halfspaces whose weights are orthogonal to the normal of the band.

**Proposition 3.1.** Let  $\mathcal{D}$  be an isotropic strongly log-concave distribution. Let  $\mathbf{v} \in \mathbb{S}^{d-1}$  be any fixed direction. Let  $\rho$  be a constant. Let  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  be a function of halfspaces of the form in Eq. (A.2), with the additional restriction that its weights  $\mathbf{w}_i \in \mathbb{S}^{d-1}$  satisfy  $\langle \mathbf{w}_i, \mathbf{v} \rangle = 0$  for all  $i$ . For some  $\alpha \in [0, 1]$ , let  $T$  denote the band  $\{x : |\langle x, \mathbf{v} \rangle| \leq \alpha\}$ . Let  $\mathcal{D}_T$  be any distribution such that  $\mathcal{D}_T$  matches moments of degree at most  $\Theta(1 - \alpha^2)$  with  $\mathcal{D}_T$  up to an additive slack of  $\alpha^{\Theta(k)}$ . Then  $\mathbb{E}_{\mathcal{D}_T}[f \circ T] - \mathbb{E}_{\mathcal{D}}[f \circ T] \leq \alpha^{\Theta(k)}$ .

We now describe some of the testers that we use. First, we need a tester that ensures that the distribution is concentrated in every single direction. More formally, the tester checks that the moments of the distribution along any direction are small.

**Proposition 3.2.** For any isotropic strongly log-concave  $\mathcal{D}$ , there exists some constants  $C_1, C_2$  and a tester  $T_1$  that takes as input  $\mathcal{D}$ , an even  $k \in \mathbb{N}$ , a parameter  $\alpha \in (0, 1)$  and runs in time  $\text{poly}(d^k; \log \frac{1}{\alpha})$ . Let  $\mathcal{D}$  denote the uniform distribution over  $\mathcal{S}$ . If  $T_1$  accepts, then for any  $\mathbf{v} \in \mathbb{S}^{d-1}$

$$\mathbb{E}_{(x,y) \sim \mathcal{D}}[(\langle \mathbf{v}, x \rangle)^k] \leq (C_1 k)^{k-2}; \quad (3.1)$$

Moreover, if  $\mathcal{S}$  is obtained by taking at least  $d^k; \log \frac{1}{\alpha}^{C_1}$  i.i.d. samples from a distribution whose  $\mathbb{R}^d$ -marginal is  $\mathcal{D}$ , the test  $T_1$  passes with probability at least  $1 - \alpha$ .

Secondly, we will use a tester that makes sure the distribution is not concentrated too close to a specific hyperplane. This is one of the properties we will need to use in order to employ the localization technique of (ABL17).

**Proposition 3.3.** For any isotropic strongly log-concave  $\mathcal{D}$ , there exist some constants  $C_2, C_3$  and a tester  $T_2$  that takes as input  $\mathcal{D}$ , a vector  $\mathbf{w} \in \mathbb{S}^{d-1}$ , parameters  $\alpha \in (0, 1)$  and runs in time  $\text{poly}(d; \log \frac{1}{\alpha})$ . Let  $\mathcal{D}$  denote the uniform distribution over  $\mathcal{S}$ . If  $T_2$  accepts, then

$$\mathbb{P}_{(x,y) \sim \mathcal{D}}[\langle \mathbf{w}, x \rangle \geq \alpha] \leq (C_2; C_3); \quad (3.2)$$

Moreover, if  $\mathcal{S}$  is obtained by taking at least  $\frac{d^{0.5}}{\alpha^{0.5}} \log \frac{1}{\alpha}$  i.i.d. samples from a distribution whose  $\mathbb{R}^d$ -marginal is  $\mathcal{D}$ , the test  $T_2$  passes with probability at least  $1 - \alpha$ .

Finally, in order to use the localization idea of (ABL17) in a manner similar to (DKTZ20b), we need to make sure that the distribution is well-behaved also within a band around to a certain hyperplane. The main property of the distribution that we establish is that functions of constantly many halfspaces have expectations very close to what they would be under our distributional assumption. As we show later in this work, having the aforementioned property allows us to derive many other properties that strongly log-concave distributions have, including many of the key properties that make the localization technique successful.

**Proposition 3.4.** For any isotropic strongly log-concave  $\mathcal{D}$  and a constant  $C_4$ , there exists a constant  $C_5$  and a tester  $T_3$  that takes as input  $\mathcal{D}$ , a vector  $\mathbf{w} \in \mathbb{S}^{d-1}$ , parameters  $\alpha \in (0, 1)$  and runs in time  $\text{poly}(d^{\Theta(\frac{1}{\alpha})}; \log \frac{1}{\alpha})$ . Let  $\mathcal{D}$  denote the uniform distribution over  $\mathcal{S}$ , let  $T$  denote the band  $\{x : |\langle x, \mathbf{w} \rangle| \leq \alpha\}$  and let  $\mathcal{F}_w$  denote the set of  $\mathcal{F}_w$ -valued functions of  $C_4$  halfspaces whose weight vectors are orthogonal to  $\mathbf{w}$ . If  $T_3$  accepts, then

$$\max_{f \in \mathcal{F}_w} \mathbb{E}_{x \sim \mathcal{D}}[f(x) \circ T] - \mathbb{E}_{(x,y) \sim \mathcal{D}}[f(x) \circ T] \leq \alpha^{C_5}; \quad (3.3)$$

$$\max_{\substack{v \geq 1 \\ S^d}} \mathbb{E}_{(x,y) \sim D} [(hv; xi)^2] \leq \frac{1}{2} \mathbb{E}_{(x,y) \sim D} [(hv; xi)^2] \quad (3.4)$$

Moreover, if  $S$  is obtained by taking at least  $\frac{1}{\epsilon} \frac{1}{d} \frac{1}{\log^{C_5}(\frac{1}{\epsilon})} \log \frac{1}{\epsilon} \frac{1}{2} \log^{C_5}(\frac{1}{\epsilon})^{C_5}$  i.i.d. samples from a distribution whose  $S^d$ -marginal is  $D$ , the test  $T_3$  passes w.p. at least  $1 - \epsilon$ .

#### 4 TESTABLY LEARNING HALFSPACES WITH MASSART NOISE

In this section we prove that we can testably learn halfspaces with Massart noise with respect to isotropic strongly log-concave distributions (see Definition A.1).

**Theorem 4.1 (Tester-Learner for Halfspaces with Massart Noise)** Let  $D_{XY}$  be a distribution over  $\mathbb{R}^d \times \{-1, 1\}$  and let  $D$  be an isotropic strongly log-concave distribution over  $\mathbb{R}^d$ . Let  $\mathcal{C}$  be the class of origin centered halfspaces in  $\mathbb{R}^d$ . Then, for any  $\epsilon < 1/2$ ,  $\delta > 0$  and  $\eta \in (0, 1)$ , there exists an algorithm (Algorithm 1) that testably learns  $\mathcal{C}$  w.r.t.  $D$  up to excess error  $\epsilon$  and error probability at most  $\delta$  in the Massart noise model with rate at most  $\frac{1}{\epsilon^2 \delta}$  using time and a number of samples from  $D_{XY}$  that are polynomial in  $d, 1/\epsilon, 1/\delta, 1/\eta$  and  $\log(1/\delta)$ .

---

##### Algorithm 1: Tester-learner for halfspaces

---

Input: Training sets  $S_1, S_2$ , parameters  $\epsilon, \delta, \eta$ ,

Output: A near-optimal weight vector  $w$ , or rejection

Run PSGD on the empirical loss  $\ell_S$  over  $S_1$  to get a list  $L$  of candidate vectors.

Test whether  $L$  contains an  $\epsilon$ -approximate stationary point of the empirical loss  $\ell_S$  over  $S_2$ .

Reject if no such  $w$  exists.

for each candidate  $w^0$  in  $L$  do

Let  $B_{w^0}(\epsilon)$  denote the band  $\{x : |hw^0; x| \leq \epsilon\}$ . Let  $F_{w^0}$  denote the class of functions of at most two halfspaces with weights orthogonal to  $w^0$ .

Let  $\epsilon^0 = \epsilon/2$ .

Run  $T_1(S_2; k=2; \epsilon^0)$  to verify that the empirical marginal is approximately isotropic.

Reject if  $T_1$  rejects.

Run  $T_2(S_2; w^0; \epsilon^0)$  to verify that  $P_S[B_{w^0}(\epsilon^0)] = \epsilon^0$ . Reject if  $T_2$  rejects.

Run  $T_3(S_2; w^0; \epsilon^0 = 6\epsilon; \delta; \eta)$  and  $T_3(S; w^0; \epsilon^0 = 2\epsilon; \delta; \eta)$  for a suitable constant to verify that the empirical distribution conditioned on  $B_{w^0}(\epsilon^0)$  and  $B_{w^0}(\epsilon^0)$  fools  $F_{w^0}$  up to  $\epsilon^0$ . Reject if  $T_3$  rejects.

Estimate the empirical error  $\ell_S(w^0)$  on  $S$ .

If all tests have accepted, output  $w^0$  with the best empirical error.

---

To show our result, we revisit the approach of (DKTZ20a) for learning halfspaces with Massart noise under well-behaved distributions. Their result is based on the idea of minimizing a surrogate loss that is non convex, but whose stationary points correspond to halfspaces with low error. They also require that their surrogate loss is sufficiently smooth, so that one can find a stationary point efficiently. While the distributional assumptions that are used to demonstrate that stationary points of the surrogate loss can be discovered efficiently are mild, the main technical lemma, which demonstrates that any stationary point suffices, requires assumptions that are not necessarily testable. We establish a label-dependent approach for testing, making use of tests that are applied during the course of our algorithm.

We consider a slightly different surrogate loss than the one used in (DKTZ20a). In particular, for  $\epsilon > 0$ , we let

$$L(w) = \mathbb{E}_{(x,y) \sim D_{XY}} \left[ y \frac{hw; xi}{\|w\|_2} \right]; \quad (4.1)$$

where  $\sigma_\epsilon : \mathbb{R} \rightarrow [0, 1]$  is a smooth approximation to the ramp function with the properties described in Proposition C.1 (see Appendix C), obtained using a piecewise polynomial of degree like the standard logistic function, our loss function has derivative exactly away from the origin (for  $|t| > \epsilon$ ). This makes the analysis of the gradient of easier, since the contribution from points lying outside a certain band is exactly

The smoothness allows us to run PSGD to obtain stationary points efficiently, and we now state the convergence lemma we need.

**Proposition 4.2 (PSGD Convergence, Lemmas 4.2 and B.2 in (DKTZ20b)).** Let  $L$  be as in Equation equation 4.1 with  $\beta \in (0; 1]$ ,  $\eta$  as described in Proposition C.1 and  $D_{XY}$  such that the marginal  $D_X$  on  $\mathbb{R}^d$  satisfies Property equation 3.1 for  $\alpha = 2$ . Then, for any  $\epsilon > 0$  and  $\beta \in (0; 1)$ , there is an algorithm whose time and sample complexity is  $O(d + \frac{\log(1/\epsilon)}{\beta})$ , which, having access to samples from  $D_{XY}$ , outputs a list  $L$  of vectors  $w \in S^{d-1}$  with  $|L| = O(\frac{d}{\beta} + \frac{\log(1/\epsilon)}{\beta})$  so that there exists  $w \in L$  with

$$\| \nabla_{w \in L} \ell(w) \|_2 \leq \epsilon; \text{ with probability at least } 1 - \epsilon.$$

In particular, the algorithm performs Stochastic Gradient Descent on Projected on  $S^{d-1}$  (PSGD).

It now suffices to show that, upon performing PSGD on  $L$ , for some appropriate choice of  $\beta$  we acquire a list of vectors that testably contain a vector which is approximately optimal. We first prove the following lemma, whose distributional assumptions are relaxed compared to the corresponding structural Lemma 3.2 of (DKTZ20a). In particular, instead of requiring the marginal distribution to be “well-behaved”, we assume that the quantities of interest (for the purposes of our proof) have expected values under the true marginal distribution that are close, up to multiplicative factors, to their expected values under some “well-behaved” (in fact, strongly log-concave) distribution. While some of the quantities of interest have values that are miniscule and estimating them up to multiplicative factors could be too costly, it turns out that the source of their vanishing scaling can be completely attributed to factors of the form  $\prod_{i,j} w_{ij}^{x_{ij}}$  (where  $\beta$  is small), which, due to standard concentration arguments, can be approximated up to multiplicative factors, given  $w \in S^{d-1}$  and  $\beta > 0$  (see Proposition 3.3). As a result, we may estimate the remaining factors up to sufficiently small additive constants (see Proposition 3.4) to get multiplicative overall closeness to the “well behaved” baseline. We defer the proof of the following Lemma to Appendix C.1.

**Lemma 4.3.** Let  $L$  be as in Equation equation 4.1 with  $\beta \in (0; 1]$ ,  $\eta$  as described in Proposition C.1, let  $w \in S^{d-1}$  and consider  $D_{XY}$  such that the marginal  $D_X$  on  $\mathbb{R}^d$  satisfies Properties equation 3.2 and equation 3.3 for  $\alpha = 2$  and accuracy  $\epsilon$ . Let  $w^* \in S^{d-1}$  define an optimum halfspace and let  $\beta < 1/2$  be an upper bound on the rate of the Massart noise. Then, there are constants  $c_1, c_2, c_3 > 0$  such that if  $\| \nabla_{w \in L} \ell(w) \|_2 < c_1(1 - \beta)$  and  $\beta < c_2$ , then

$$\ell(w; w^*) \leq \frac{c_3}{1 - \beta} \quad \text{or} \quad \ell(w; w^*) \leq \frac{c_3}{1 - \beta}$$

Combining Proposition 4.2 and Lemma 4.3, we get that for any choice of the parameters  $\beta \in (0; 1]$ , by running PSGD on  $L$ , we can construct a list of vectors of polynomial size (in all relevant parameters) that testably contains a vector that is close to the optimum weight vector. In order to link the zero-one loss to the angular similarity between a weight vector and the optimum vector, we use the following Proposition (for the proof, see Appendix C.2).

**Proposition 4.4.** Let  $D_{XY}$  be a distribution over  $\mathbb{R}^d$  for  $\beta \in (0; 1]$ ,  $w \in \arg \min_{w \in S^{d-1}} P_{D_{XY}} [y \neq \text{sign}(hw; x_i)]$  and  $w \in S^{d-1}$ . Then, for any  $\epsilon \in (0; 1]$ ,  $\beta \in (0; 1/4]$ , if the marginal  $D_X$  on  $\mathbb{R}^d$  satisfies Property equation 3.1 for  $\alpha_1 > 0$  and some even  $k \geq N$  and Property equation 3.2 with set to  $(C_1 k)^{\frac{k}{2(k+1)}} (\tan \beta)^{\frac{k}{k+1}}$ , then, there exists a constant  $\epsilon_0 > 0$  such that the following is true.

$$P_{D_{XY}} [y \neq \text{sign}(hw; x_i)] \leq \epsilon + c \beta^{1-2\alpha_1} (1 - \beta)^{\frac{1}{k+1}};$$

**Proof of Theorem 4.1.** Throughout the proof we consider  $\beta$  to be a sufficiently small polynomial in all the relevant parameters. Each of the failure events will have probability at least their number will be polynomial in all the relevant parameters, so by the union bound, we may pick that the probability of failure is at most

The algorithm we run is Algorithm 1, with appropriate selection of parameters and given samples  $S_1, S_2$ , each of which are sufficiently large sets of independent samples from the true unknown distribution  $D_{XY}$ . For some  $\beta \in (0; 1]$  to be defined later, we run PSGD on the empirical loss over  $S_1$  as described in Proposition 4.2 with  $\beta = c_1(1 - \beta) = 4$ , where  $c_1$  is given by Lemma 4.3.

Figure 1: Critical regions in the proofs of main structural lemmas (Lemmas 4.3, 5.2). We analyze the contributions of the regions labeled  $A_1, A_2$  to the quantities  $A_1, A_2$  in the proofs. Specifically, the regions  $A_1$  (which have height  $= 3$  so that the value of  $f^0(x_w)$  for any  $x$  in these regions is exactly  $1=$ , by Proposition C.1) form a subset of the region  $C$  and their probability mass under  $D_{f,x}$  is (up to a multiplicative factor) a lower bound on the quantity (see Eq equation C.3). Similarly, the region  $A_2$  is a subset of the intersection  $C \cap B$  with the band of height  $=$ , and has probability mass that is (up to a multiplicative factor) an upper bound on the quantity (see Eq equation C.4).

By Proposition 4.2, we get a list of vectors  $w \in \mathbb{S}^{d-1}$  with  $|L_j| = \text{poly}(d; 1=)$  such that there is  $w \in L$  with  $\|w\|_2 < \frac{1}{2}c_1(1 - 2)$  under the true distribution, if the marginal is isotropic.

Having acquired the list using samples  $S_1$ , we use the independent samples  $S_2$  to test whether  $L$  contains an approximately stationary point of the empirical loss  $\mathcal{L}_2$ . If this is not the case, then we may safely reject: for large enough  $|S_2|$ , if the distribution is indeed isotropic strongly logconcave, there is an approximate stationary of the population loss  $\mathcal{L}$  and if  $|S_2|$  is large enough, the gradient of the empirical loss  $\mathcal{L}_2$  will be close to the gradient of the population loss on each of the elements of  $L$ , due to appropriate concentration bounds for log-concave distributions as well as the fact that the elements of  $L$  are independent from  $S_2$ . For the following, let  $w$  be a point such that  $\|w\|_2 < c_1(1 - 2)$  under the empirical distribution over  $S_2$ .

In Lemma 4.3 and Proposition 4.4 we have identified certain properties of the marginal distribution that are sufficient for our purposes, given that  $L$  contains an approximately stationary point of the empirical (surrogate) loss  $\mathcal{L}_2$ . Our testers  $T_1, T_2, T_3$  verify that these properties hold for the empirical marginal over our samples  $S_2$ , and it will be convenient to analyze the optimality of our algorithm purely over  $S_2$ . In particular, we will need to require that  $|S_2|$  is sufficiently large, so that when the true marginal is indeed the target, our testers succeed with high probability (for the corresponding sample complexity, see Propositions 3.2, 3.3 and 3.4). Moreover, by standard generalization theory, since the VC dimension of halfspaces is  $O(d)$  and for us  $|S_2|$  is a large  $\text{poly}(d; 1=)$ , both the error of our final output and the optimal error over  $S_2$  will be close to that over  $D_{X,Y}$ . So in what follows, we will abuse notation and refer to the uniform distribution over  $S_2$  as  $D_{X,Y}$  and the optimal error over  $S_2$  simply as  $\text{opt}$ .

We proceed with some basic tests. Throughout the rest of the algorithm, whenever a tester fails, we reject, otherwise we proceed. First, we run tester  $T_1$  with inputs  $(w; = 2; \epsilon)$  and  $(w; = 6; \epsilon)$  (Proposition 3.3) and  $T_3$  with inputs  $(w; = 2; c_2; \epsilon)$  and with  $(w; = 6; c_2; \epsilon)$  (Proposition 3.4  $c_2$  as defined in Lemma 4.3). This ensures that for the approximate stationary point  $w$  of the  $L$ , the probability within the band  $B_w(= 2) = \{x : |x_j| = 2g\}$  (and similarly for  $B_w(= 6)$ ) and moreover that our marginal conditioned on each of the bands fools (up to an additive constant) functions of halfspaces with weights orthogonal to  $w$ . As a result, we may apply Lemma 4.3 to and form a list of 2 vectors  $f; w; w^0$  which contains some  $w^0$  with  $\|w^0\|_2 = c_2 = (1 - 2)$  (where  $c_3$  is as defined in Lemma 4.3).



We run  $T_1$  (Proposition 3.2) with  $k = 2$  to verify that the marginals are approximately isotropic and we use  $T_2$  once again, with appropriate parameters for  $\mathcal{D}$  and its negation, to apply Proposition 4.4 and get that  $w; w^0$  contains a vector  $w^0$  with

$$\mathbb{P}_{D_{XY}} [y \notin \text{sign}(hw^0; xi)] \leq \text{opt} + c^{-2=3}; \text{ where } (w^0; w) := c_2 = \frac{1}{1-2}$$

By picking  $\epsilon = (1-2)^{3=2}$  we have  $\mathbb{P}_{D_{XY}} [y \notin \text{sign}(hw^0; xi)] \leq \text{opt} + \epsilon$

However, we do not know which of the weight vectors  $w; w^0$  is the one guaranteed to achieve small error. In order to discover this vector, we estimate the probability of error of each of the corresponding halfspaces (which can be done efficiently, due to Hoeffding's bound) and pick the one with the smallest error. This final step does not require any distributional assumptions and we do not need to perform any further tests.  $\square$

## 5 TESTABLY LEARNING HALFSPACES IN THE AGNOSTIC SETTING

In this section, we provide our result on efficiently and testably learning halfspaces in the agnostic setting with respect to isotropic strongly log-concave target marginals. We defer the proofs to Appendix D. The algorithm we use is once more Algorithm 1, but we call it multiple times for different choices of the parameter, reject if any call rejects and output the vector that achieved the minimum empirical error overall, otherwise. Also, the tester is called for a general  $k$  (not necessarily  $k = 2$ ).

**Theorem 5.1 (Efficient Tester-Learner for Halfspaces in the Agnostic Setting)** Let  $D_{XY}$  be a distribution over  $\mathbb{R}^d \times \{-1, 1\}$  and let  $D$  be a strongly log-concave distribution over  $\mathbb{R}^d$  (Definition A.1). Let  $\mathcal{C}$  be the class of origin centered halfspaces. Then, for any even  $k \geq 2$ , any  $\epsilon > 0$  and  $\delta \in (0, 1)$ , there exists an algorithm that agnostically testably learns w.r.t.  $D$  up to error  $O(k^{1=2} \text{opt}^{1+\frac{1}{k+1}}) + \epsilon$ , where  $\text{opt} = \min_{w \in S^{d-1}} \mathbb{P}_{D_{XY}} [y \notin \text{sign}(hw; xi)]$ , and error probability at most  $\delta$ , using time and a number of samples from  $D_{XY}$  that are polynomial in  $d^{O(k)}; (1/\epsilon)^{O(k)}$  and  $(\log(1/\delta))^{O(k)}$ . In particular, by picking some appropriate  $\epsilon = \log^2 d$ , we obtain error  $O(\text{opt}) + \epsilon$  in quasipolynomial time and sample complexity,  $\text{poly}(2^{\text{polylog } d}; (1/\delta)^{\text{polylog } d})$ .

To prove Theorem 5.1, we may follow a similar approach as the one we used for the case of Massart noise. However, in this case, the main structural lemma regarding the quality of the stationary points involves an additional requirement about the parameter  $\epsilon$ . In particular,  $\epsilon$  cannot be arbitrarily small with respect to the error of the optimum halfspace, because, in this case, there is no upper bound on the amount of noise that any specific point might be associated with. As a result, picking  $\epsilon$  to be arbitrarily small would imply that our algorithm only considers points that lie within a region that has arbitrarily small probability and can hence be completely corrupted with the adversarial  $\text{opt}$  budget. On the other hand, the polynomial slackness that the testability requirement introduces (through Proposition 4.4) between the error we achieve and the angular distance guarantee we can get via finding a stationary point  $w^0$  (which is now coupled with  $\text{opt}$ ), appears to the exponent of the guarantee we achieve in Theorem 5.1.

**Lemma 5.2.** Let  $L$  be as in Equation equation 4.1 with  $\delta \in (0, 1)$ ,  $\epsilon$  as described in Proposition C.1, let  $w \in S^{d-1}$  and consider  $D_{XY}$  such that the marginal  $D_X$  on  $\mathbb{R}^d$  satisfies Properties equation 3.2, equation 3.3 and equation 3.4 for with  $C_4 = 2$  and accuracy parameter. Let  $\text{opt}$  be the minimum error achieved by some origin centered halfspace and let  $w \in S^{d-1}$  be a corresponding vector. Then, there are constants  $c_2; c_3; c_4 > 0$  such that if  $\text{opt} \leq c_1$ ,  $\text{kr}_w L(w) k_2 < c_2$ , and  $c_3$  then either  $\mathbb{P}_{D_{XY}} [y \notin \text{sign}(hw; xi)] \leq c_4$  or  $\mathbb{P}_{D_{XY}} [y \notin \text{sign}(hw; xi)] \geq c_4$ :

We obtain our main result for Gaussian target marginals by refining Proposition 4.4 for the specific case when the target marginal distribution is the standard multivariate Gaussian distribution. The algorithm for the Gaussian case is similar to the one of Theorem 5.1, but it runs different tests for the improved version (see Proposition D.1) of Proposition 4.4.

**Theorem 5.3.** In Theorem 5.1, if  $D$  is the standard Gaussian in  $d$  dimensions, we obtain error  $O(\text{opt}) + \epsilon$  in polynomial time and sample complexity,  $\text{poly}(d; 1/\delta; \log(1/\delta))$ .

## ACKNOWLEDGMENTS

We wish to thank the anonymous reviewers of ICLR 2024 for their constructive feedback. Aravind Gollakota was at UT Austin while this work was done, supported by NSF award AF-1909204 and the NSF AI Institute for Foundations of Machine Learning (IFML). Adam R. Klivans was supported by NSF award AF-1909204 and the NSF AI Institute for Foundations of Machine Learning (IFML). Konstantinos Stavropoulos was supported by NSF award AF-1909204, the NSF AI Institute for Foundations of Machine Learning (IFML), and by scholarships from Bodossaki Foundation and Leventis Foundation. Arsen Vasilyan was supported in part by NSF awards CCF-2006664, DMS-2022448, CCF-1565235, CCF-1955217, CCF-2310818, Big George Fellowship and Fin-tech@CSAIL. Part of this work was done while Arsen Vasilyan was visiting UT Austin.

## REFERENCES

- [ABHU15] Pranjal Awasthi, Maria-Florina Balcan, Nika Haghtalab, and Ruth Uerner. Efficient learning of linear separators under bounded noise. *Conference on Learning Theory* pages 167–190. PMLR, 2015.
- [ABHZ16] Pranjal Awasthi, Maria-Florina Balcan, Nika Haghtalab, and Hongyang Zhang. Learning and 1-bit compressed sensing under asymmetric noise. *Conference on Learning Theory* pages 152–192. PMLR, 2016.
- [ABL17] Pranjal Awasthi, Maria Florina Balcan, and Philip M Long. The power of localization for efficiently learning linear separators with noise. *Journal of the ACM (JACM)* 63(6):1–27, 2017.
- [BH21] Maria-Florina Balcan and Nika Haghtalab. Noise in classification. *Beyond the Worst-Case Analysis of Algorithms* page 361, 2021.
- [BZ17] Maria-Florina F Balcan and Hongyang Zhang. Sample and computationally efficient learning algorithms under  $s$ -concave distributions. *Advances in Neural Information Processing Systems* 30, 2017.
- [Dan15] Amit Daniely. A ptas for agnostically learning halfspaces. *Conference on Learning Theory* pages 484–502. PMLR, 2015.
- [DGT19] Ilias Diakonikolas, Themis Gouleakis, and Christos Tzamos. Distribution-independent pac learning of halfspaces with massart noise. *Advances in Neural Information Processing Systems* 32, 2019.
- [DK22] Ilias Diakonikolas and Daniel Kane. Near-optimal statistical query hardness of learning halfspaces with massart noise. *Conference on Learning Theory* pages 4258–4282. PMLR, 2022.
- [DKK<sup>+</sup>22] Ilias Diakonikolas, Daniel M Kane, Vasilis Kontonis, Christos Tzamos, and Nikos Zari s. Learning general halfspaces with general massart noise under the gaussian distribution. In *Proceedings of the 54th Annual ACM SIGACT Symposium on Theory of Computing* pages 874–885, 2022.
- [DKK<sup>+</sup>23] Ilias Diakonikolas, Daniel M Kane, Vasilis Kontonis, Sihan Liu, and Nikos Zari s. Efficient testable learning of halfspaces with adversarial label noise. *arXiv preprint arXiv:2303.05485* 2023.
- [DKMR22] Ilias Diakonikolas, Daniel Kane, Pasin Manurangsi, and Lisheng Ren. Cryptographic hardness of learning halfspaces with massart noise. *Advances in Neural Information Processing Systems* 35, 2022.
- [DKPZ21] Ilias Diakonikolas, Daniel M Kane, Thanasis Pittas, and Nikos Zari s. The optimality of polynomial regression for agnostic learning under gaussian marginals in the sq model. In *Conference on Learning Theory* pages 1552–1584. PMLR, 2021.

- [DKS18] Ilias Diakonikolas, Daniel M Kane, and Alistair Stewart. Learning geometric concepts with nasty noise. In *Proceedings of the 50th Annual ACM SIGACT Symposium on Theory of Computing* pages 1061–1073, 2018.
- [DKT21] Ilias Diakonikolas, Daniel Kane, and Christos Tzamos. Forster decomposition and learning halfspaces with noise. *Advances in Neural Information Processing Systems* 34:7732–7744, 2021.
- [DKTZ20a] Ilias Diakonikolas, Vasilis Kontonis, Christos Tzamos, and Nikos Zari s. Learning halfspaces with massart noise under structured distribution. *Conference on Learning Theory* pages 1486–1513. PMLR, 2020.
- [DKTZ20b] Ilias Diakonikolas, Vasilis Kontonis, Christos Tzamos, and Nikos Zari s. Non-convex sgd learns halfspaces with adversarial label noise. *Advances in Neural Information Processing Systems* 33:18540–18549, 2020.
- [DKTZ22] Ilias Diakonikolas, Vasilis Kontonis, Christos Tzamos, and Nikos Zari s. Learning general halfspaces with adversarial label noise via online gradient descent. *International Conference on Machine Learning* pages 5118–5141. PMLR, 2022.
- [DKZ20] Ilias Diakonikolas, Daniel Kane, and Nikos Zari s. Near-optimal sq lower bounds for agnostically learning halfspaces and relus under gaussian margin. *Advances in Neural Information Processing Systems* 33:13586–13596, 2020.
- [DTK22] Ilias Diakonikolas, Christos Tzamos, and Daniel M Kane. A strongly polynomial algorithm for approximate forster transforms and its application to halfspace learning. *arXiv preprint arXiv:2212.03008*, 2022.
- [GGK20] Surbhi Goel, Aravind Gollakota, and Adam Klivans. Statistical-query lower bounds via functional gradients. *Advances in Neural Information Processing Systems* 33:2147–2158, 2020.
- [GKK23] Aravind Gollakota, Adam R Klivans, and Pravesh K Kothari. A moment-matching approach to testable learning and a new characterization of rademacher complexity. *Proceedings of the 54th annual ACM Symposium on Theory of Computing*, 2023. To appear.
- [GKSV23] Aravind Gollakota, Adam R Klivans, Konstantinos Stavropoulos, and Arsen Vasilyan. Tester-learners for halfspaces: Universal algorithm. *arXiv preprint arXiv:2305.11765*, 2023.
- [KKMS08] Adam Tauman Kalai, Adam R Klivans, Yishay Mansour, and Rocco A Servedio. Agnostically learning halfspaces. *SIAM Journal on Computing* 37(6):1777–1805, 2008.
- [KLS09] Adam R Klivans, Philip M Long, and Rocco A Servedio. Learning halfspaces with malicious noise. *Journal of Machine Learning Research* 10(12), 2009.
- [RV23] Ronitt Rubinfeld and Arsen Vasilyan. Testing distributional assumptions of learning algorithms. *Proceedings of the 54th annual ACM Symposium on Theory of Computing* 2023. To appear.
- [SW14] Adrien Saumard and Jon A Wellner. Log-concavity and strong log-concavity: a review. *Statistics surveys* 3:45, 2014.
- [YZ17] Songbai Yan and Chicheng Zhang. Revisiting perceptron: Efficient and label-optimal learning of halfspaces. *Advances in Neural Information Processing Systems* 30:2017, 2017.
- [Zha18] Chicheng Zhang. Efficient active learning of sparse halfspaces. *Conference on Learning Theory* pages 1856–1880. PMLR, 2018.
- [ZL21] Chicheng Zhang and Yinan Li. Improved algorithms for efficient active learning halfspaces with massart and tsybakov noise. *Conference on Learning Theory* pages 4526–4527. PMLR, 2021.

[ZSA20] Chicheng Zhang, Jie Shen, and Pranjali Awasthi. Efficient active learning of sparse halfspaces with arbitrary bounded noise. *Advances in Neural Information Processing Systems*, 33:7184–7197, 2020.

## A STRONGLY LOG-CONCAVE DISTRIBUTIONS

We also formally define the class of strongly log-concave distributions, which is the class that our target marginal  $D$  is allowed to belong to, and collect some useful properties of such distributions. We will state the definition for isotropic  $D$  (i.e. with mean 0 and covariance  $I$ ) for simplicity.

**Definition A.1** (Strongly log-concave distribution, see e.g. (SW14, Def 2.5)). We say an isotropic distribution  $D$  on  $\mathbb{R}^d$  is strongly log-concave if the logarithm of its density is a strongly concave function. Equivalently, it can be written as

$$q(x) = r(x) \exp(-\lambda \|x\|^2) \quad (\text{A.1})$$

for some log-concave function  $r$  and some constant  $\lambda > 0$ , where  $\exp(-\lambda \|x\|^2)$  denotes the density of the spherical Gaussian  $\mathcal{N}(0; \lambda^{-1}I)$ .

**Proposition A.2** (see e.g. (SW14)). Let  $D$  be an isotropic strongly log-concave distribution with density  $q$ .

- (a) Any orthogonal projection of  $D$  onto a subspace is also strongly log-concave.
- (b) There exist constants  $c, R$  such that  $q(x) \leq U$  for all  $x$ , and  $q(x) \geq U$  for all  $\|x\| \leq R$ .
- (c) There exist constants  $c^0$  and  $c$  such that  $q(x) \geq U^0 \exp(-\lambda \|x\|^2)$  for all  $x$ .
- (d) There exist constants  $K_1, K_2$  such that for any  $\alpha \in [0, 1]$  and any  $v \in \mathbb{S}^{d-1}$ ,  $\mathbb{P}[|v \cdot x| \geq \alpha] \leq (K_1 + K_2 \alpha^{-2})$ .
- (e) There exists a constant  $K_3$  such that for any  $k \in \mathbb{N}$ ,  $\mathbb{E}[|v \cdot x|^k] \leq (K_3 k)^{k-2}$ .
- (f) Let  $\bar{Q} = (i_1, \dots, i_d) \in \mathbb{Z}_{\geq 0}^d$  be a multi-index with total degree  $j = \sum_{i=1}^d i_i = k$ , and let  $x = \prod_{i=1}^d x_i^{i_i}$ . There exists a constant  $K_4$  such that for any such  $\bar{Q}$ ,  $\mathbb{E}[x^{\bar{Q}}] \leq (K_4 k)^{k-2}$ .

For (a), see e.g. (SW14, Thm 3.7). The other properties follow readily from Eq. (A.1), which allows us to treat the density as subgaussian.

A key structural fact that we will need about strongly log-concave distributions is that approximately matching moments of degree at most  $k$  with such a  $D$  is sufficient to fool any function of a constant number of halfspaces up to an additive

**Proposition A.3** (Variant of (GKK23, Thm 5.6)). Let  $p$  be a fixed constant, and let  $\mathcal{F}$  be the class of all functions of  $p$  halfspaces mapping  $\mathbb{R}^d$  to  $[-1, 1]$  of the form

$$f(x) = g(\text{sign}(v^1 \cdot x_1 + \beta_1); \dots; \text{sign}(v^p \cdot x_1 + \beta_p)) \quad (\text{A.2})$$

for some  $g: \{-1, 1\}^p \rightarrow [-1, 1]$  and weights  $v^i \in \mathbb{S}^{d-1}$ . Let  $D$  be any target marginal such that for every  $i$ , the projection  $v^i \cdot x_i$  has subgaussian tails and is anticoncentrated:  $\mathbb{P}[|v^i \cdot x_i| > t] \leq \exp(-c t^2)$ , and (b) for any interval  $[a, b]$ ,  $\mathbb{P}[v^i \cdot x_i \in [a, b]] \geq \frac{1}{p} |b - a|$ . Let  $D$  be any distribution such that for all monomials  $s = \prod_{i=1}^d x_i^{i_i}$  of total degree  $j = \sum_{i=1}^d i_i \leq k$ ,

$$\mathbb{E}_D[x^s] = \mathbb{E}_D[x^s] + \frac{c_j}{d^j} \frac{1}{k^j}$$

for some sufficiently small constant  $c_j$  (in particular, it suffices to have  $c_j = \Theta(k^{-j})$  moment closeness for every  $j$ ). Then

$$\max_{f \in \mathcal{F}} \mathbb{E}_D[f] - \mathbb{E}_D[f] \leq \frac{1}{p^k}$$

Note that this is a variant of the original statement of (GKK23, Thm 5.6), which requires that the 1D projection of  $D$  along any direction satisfy suitable concentration and anticoncentration. Indeed, an inspection of their proof reveals that it suffices to verify these properties for projections only along the directions  $v^i$ ,  $i \in [p]$  as opposed to all directions. This is because to fool a function of the form above, their proof only analyzes the projected distribution  $(v^1 \cdot x_1; \dots; v^p \cdot x_1)$  on  $\mathbb{R}^p$ , and requires only concentration and anticoncentration for each individual projection  $v^i \cdot x_i$ .

## B PROOFS FOR SECTION 3

### B.1 PROOF OF PROPOSITION 3.1

Our plan is to apply Proposition A.3. To do so, we must verify that  $D_{jT}$  satisfies the assumptions required. In particular, it suffices to verify that the 1D projection along any direction orthogonal to  $w$  has subgaussian tails and is anticoncentrated. Let  $S^d \ni v$  be any direction that is orthogonal to  $w$ . By Proposition A.2(d), we may assume that  $\|T\|_F \leq 1$ .

To verify subgaussian tails, we must show that for any  $t > 0$ ,  $\mathbb{P}_{D_{jT}}[\langle hv; x \rangle > t] \leq \exp(-Ct^2)$  for some constant  $C$ . The main fact we use is Proposition A.2(c), i.e. that any strongly log-concave density is pointwise upper bounded by a Gaussian density times a constant. Write

$$\mathbb{P}_{D_{jT}}[\langle hv; x \rangle > t] = \frac{\mathbb{P}_D[\langle hv; x \rangle > t \text{ and } \|x\|_2 \leq \frac{1}{\epsilon}]}{\mathbb{P}_D[\|x\|_2 \leq \frac{1}{\epsilon}]}$$

The claim now follows from the fact that the numerator is upper bounded by a constant times the corresponding probability under a Gaussian density, which is at most  $\exp(-Ct^2)$  for some constant  $C$ , and that the denominator is  $\geq \epsilon^d$ .

To check anticoncentration, for any interval  $[a, b]$ , write

$$\mathbb{P}_{D_{jT}}[\langle hv; x \rangle \in [a, b]] = \frac{\mathbb{P}_D[\langle hv; x \rangle \in [a, b] \text{ and } \|x\|_2 \leq \frac{1}{\epsilon}]}{\mathbb{P}_D[\|x\|_2 \leq \frac{1}{\epsilon}]}$$

After projecting onto  $\text{span}(v; w)$  (an operation that preserves log-concavity), the numerator is the probability mass under a rectangle with side lengths  $a_j$  and  $2\epsilon$ , which is at most  $O(\epsilon^d)$  as by Proposition A.2(b) the density is pointwise upper bounded by a constant. The claim follows since the denominator is  $\geq \epsilon^d$ .

Now we are ready to apply Proposition A.3. We see that  $D_{jT}$  matches moments of degree at most  $k$  with  $D_{jT}$  up to an additive slack of  $O(\epsilon^{d-k})$ , then  $\mathbb{E}_D[f_j(T)] - \mathbb{E}_{D_{jT}}[f_j(T)] \leq O(\epsilon^{d-k})$ . Rewriting in terms of  $\mathbb{E}_D$  gives the theorem.

### B.2 PROOF OF PROPOSITION 3.2

The tester  $\bar{T}_1$  does the following:

1. For all  $2 \leq j \leq k$ :
  - (a) Compute the corresponding moment  $\mathbb{E}_{(x,y) \sim D} [x^j] := \frac{1}{j!} \mathbb{P}_{x \sim D} [x^j]$ .
  - (b) If  $\mathbb{E}_{(x,y) \sim D} [x^j] - \mathbb{E}_x [x^j] > \frac{1}{d^k}$  then reject.
2. If all the checks above passed, accept.

First, we claim that for some absolute constant  $C$ , if the tester above accepts, we have  $\mathbb{E}_{(x,y) \sim D} [(hv; xi)^k] \leq (C_1 k)^{k-2}$  for any  $v \in S^{d-1}$ . To show this, we first recall that by Proposition A.2(e) it is the case that  $\mathbb{E}_{(x,y) \sim D} [(hv; xi)^k] \leq (K_3 k)^{k-2}$ . But we have

$$\begin{aligned} \mathbb{E}_{(x,y) \sim D} [(hv; xi)^k] &= \mathbb{E}_{(x,y) \sim D} \left[ \sum_{j=0}^k \binom{k}{j} (hv; xi)^j \mathbb{E}_x [x^{k-j}] \right] \\ &\leq d^k \max_{j=0, \dots, k} \mathbb{E}_{(x,y) \sim D} [x^j] \mathbb{E}_x [x^{k-j}] \leq 1 \end{aligned}$$

Together with the bound  $\mathbb{E}_{(x,y) \sim D} [(hv; xi)^k] \leq (K_3 k)^{k-2}$ , the above implies that  $\mathbb{E}_{(x,y) \sim D} [(hv; xi)^k] \leq (C_1 k)^{k-2}$  for some constant  $C_1$ .

Now, we need to show that if the elements  $S$  are chosen i.i.d. from  $D$ , and  $j \leq d^k$ ;  $\log \frac{1}{\epsilon} \leq C_1$  then the tester above accepts with probability at least  $1 - \epsilon$ . Consider any speci c



Now, we need to show that if the elements  $S_j$  are chosen i.i.d. from  $\mathcal{D}$ , and  $j \in \mathcal{S}$  then the tester above accepts with probability at least  $\frac{1}{K_1}$ . Consider any specific multi-index  $z \in \mathcal{Z}^d_0$  with  $|j| = k$ . Now, by Proposition A.2(f) we have for any positive integer  $k$  the following:

$$\mathbb{E}_{x \in \mathcal{D}} \left[ \sum_{i=1}^k (x_i)^2 \right] \leq (K_4 k)^{k-2}$$

But by Proposition A.2(d) we have that  $\mathbb{P}_{x \in \mathcal{D}} [x \in \mathcal{T}] = \frac{1}{K_1}$ . Therefore, the density of the distribution  $\mathbb{D}_{\mathcal{T}}$  (which is defined as the distribution one obtains by taking and conditioning on  $\mathcal{T}$ ) is upper bounded by the product of the density of the distribution on  $\mathcal{D}$  and  $\frac{1}{K_1}$ . This allows us to bound

$$\mathbb{E}_{x \in \mathcal{D}} \left[ \sum_{i=1}^k (x_i)^2 \mid x \in \mathcal{T} \right] \leq \frac{1}{K_1} \mathbb{E}_{x \in \mathcal{D}} \left[ \sum_{i=1}^k (x_i)^2 \right] \leq \frac{(K_4 k)^{k-2}}{K_1}$$

This implies that

$$\begin{aligned} & \mathbb{E}_{x \in \mathcal{D}} \left[ \sum_{z \in \mathcal{Z}^d_0} \mathbb{E}_{z \in \mathcal{Z}^d_0} [z_j^2 \mid z \in \mathcal{T}]^{2 \log(1/\epsilon)} \right] \\ & \leq \sum_{z \in \mathcal{Z}^d_0} \mathbb{E}_{x \in \mathcal{D}} \left[ \sum_{i=1}^k (x_i)^2 \mid x \in \mathcal{T} \right]^{2 \log(1/\epsilon)} \\ & \leq \frac{1}{(K_1)^{2 \log(1/\epsilon)}} \sum_{z \in \mathcal{Z}^d_0} (K_4 k)^{k-2} (K_4 k)^{k(2 \log(1/\epsilon))} \\ & \leq \frac{1}{(K_1)^{2 \log(1/\epsilon)}} (2 K_4 \log(1/\epsilon) k)^{\log(1/\epsilon) k} \end{aligned}$$

This, together with Markov's inequality implies that

$$\mathbb{P} \left[ \frac{1}{|\mathcal{S}|} \sum_{x \in \mathcal{S}} \mathbb{E}_{x \in \mathcal{D}} [x_j^2] > \frac{1}{d^k} d^{\Theta(k)} \frac{(3 K_4 k \log(1/\epsilon))^{k-2}}{K_1 |\mathcal{S}_{\mathcal{T}}|} \right] \leq \frac{1}{|\mathcal{S}|}$$

Now, recall that the tester  $\mathcal{T}_2$  in step (1) accepted, we have  $|\mathcal{S}_{\mathcal{T}}| \geq \frac{1}{C_2} |\mathcal{S}|$ . Since  $\mathcal{S}$  is obtained by taking at least  $|\mathcal{S}| \geq \frac{1}{d} \frac{1}{\epsilon} d^{\frac{1}{2} \log C_5} \log \frac{1}{\epsilon} \frac{1}{2} \log C_5$ , for sufficiently large  $C_5$  we see that the expression above is upper-bounded by  $\frac{1}{|\mathcal{S}|}$ . Taking a union bound over all  $z \in \mathcal{Z}^d_0$  with  $|j| = k$ , we see that with probability at least  $\frac{1}{2}$  the tester  $\mathcal{T}_3$  accepts, finishing the proof.

## C PROOFS FROM SECTION 4

We first present the following Proposition, which ensures that we can form a loss function with certain desired properties.

**Proposition C.1.** There are constants  $c, c' > 0$ , such that for any  $\epsilon > 0$ , there exists a continuously differentiable function  $\eta : \mathbb{R} \rightarrow [0, 1]$  with the following properties.

1. For any  $t \in [c, c']$ ,  $\eta(t) = \frac{1}{2} + \frac{t}{c}$ .
2. For any  $t \geq 2$ ,  $\eta(t) = 1$  and for any  $t \leq -2$ ,  $\eta(t) = 0$ .
3. For any  $t \in \mathbb{R}$ ,  $\eta(t) \in [0, c']$ ,  $\eta'(t) = \eta'(-t)$  and  $|\eta'(t)| \leq c'^2$ .

**Proof.** We define  $\eta$  as follows.

$$\eta(t) = \begin{cases} \frac{1}{2} + \frac{t}{c} & \text{if } |t| \leq c \\ 1 & \text{if } t > 2 \\ 0 & \text{if } t < -2 \\ \eta'(t); t \in [c, 2] \\ \eta'(-t); t \in [-2, -c] \end{cases}$$



Figure 2: The function  $\eta$  used to smoothly approximate the ramp.

for some appropriate functions  $\eta$ . It is sufficient that we pick  $\eta^+$  satisfying the following conditions (then  $\eta$  would be defined symmetrically, i.e.,  $\eta(t) = 1 - \eta^+(t)$ ).

- $\eta^+(t) = 1$  and  $\eta^+(t) = 0$ .
- $\eta^+(t) = 2 - 3t$  and  $\eta^+(t) = 1 - t$ .
- $\eta^+$  is defined and bounded, except, possibly at  $t = 0$  and/or  $t = 2$ .

We therefore need to satisfy four equations for  $\eta^+$ . So we set  $\eta^+$  to be a degree 3 polynomial:  $\eta^+(t) = a_1 t^3 + a_2 t^2 + a_3 t + a_4$ . Whenever  $\epsilon > 0$ , the system has a unique solution that satisfies the desired inequalities. In particular, we may solve the equation to get  $a_1 = 3$ ;  $a_2 = 15 - 2\epsilon$ ;  $a_3 = 3 - 4\epsilon$  and  $a_4 = 5 - 8\epsilon$ . For the resulting function (see Figure 2 below and Figure 4 in the appendix) we have that there are constants  $c_1, c_2 > 0$  such that  $\eta^+(t) \geq c_1$  and  $|\eta^+(t)| \leq c_2$  for any  $t \in [0, 2]$ .  $\square$

### C.1 PROOF OF LEMMA 4.3

We will prove the contrapositive of the claim, namely, that there are constants  $c_3 > 0$  such that if  $\|w; w\|_2 \leq c_3$ , and  $\|w; w\|_2 > \frac{c_3}{1 - 2\epsilon}$ , then  $\|r_w L(w)\|_2 \leq c_1(1 - 2\epsilon)$ .

Consider the case where  $\|w; w\|_2 \leq 2$  (otherwise, perform the same argument for  $w$ ). Let  $v$  be a unit vector orthogonal to  $w$  that can be expressed as a linear combination of  $w$  and  $w$  for which  $\langle v; w \rangle = 0$ . Then  $\{v; w\}$  is an orthonormal basis for  $\mathcal{V} = \text{span}(w; w)$ . For any vector  $x \in \mathbb{R}^d$ , we will use the following notation  $x_w = \langle w; x \rangle$ ,  $x_v = \langle v; x \rangle$ . It follows that  $\text{proj}_{\mathcal{V}}(x) = x_w w + x_v v$ , where  $\text{proj}_{\mathcal{V}}$  is the operator that orthogonally projects vectors on  $\mathcal{V}$ .

Using the fact that  $\|w; x\|_2 = \|x\|_2 \langle w; x \rangle / \|w\|_2 = \|x\|_2 x_w$  for any  $w \in \mathbb{S}^{d-1}$ , the interchangeability of the gradient and expectation operators and the fact that  $\eta$  is an even function we get that

$$\|r_w L(w)\|_2 = \mathbb{E}_h \left[ \eta^0(\langle w; x \rangle) \langle y - \langle x; w \rangle w \rangle \right]$$

Since the projection operator  $\text{proj}_{\mathcal{V}}$  is a contraction, we have  $\|r_w L(w)\|_2 \leq \|\text{proj}_{\mathcal{V}} r_w L(w)\|_2$ , and we can therefore restrict our attention to a simpler, two dimensional problem. In particular, since  $\text{proj}_{\mathcal{V}}(x) = x_w w + x_v v$ , we get

$$\begin{aligned} \|\text{proj}_{\mathcal{V}} r_w L(w)\|_2 &= \mathbb{E}_h \left[ \eta^0(\langle x; w \rangle) \langle y - x_w w - x_v v \rangle \right] = \mathbb{E}_h \left[ \eta^0(\langle x; w \rangle) \langle y - x_w w \rangle \right] \\ &= \mathbb{E}_h \left[ \eta^0(\langle x; w \rangle) \text{sign}(\langle w; x \rangle) (1 - 2\epsilon) \langle y - \text{sign}(\langle w; x \rangle) x \rangle \right] \end{aligned}$$

Figure 3: Critical regions in the proofs of main structural lemmas (Lemmas 4.3, 5.2). We analyze the contributions of the regions labeled  $A_1, A_2$  to the quantities  $A_1, A_2$  in the proofs. Specifically, the regions  $A_1$  (which have height  $\leq 3$  so that the value of  $\phi^0(x_w)$  for any  $x$  in these regions is exactly  $1 - \frac{1}{2}$ , by Proposition C.1) form a subset of the region  $G$  and their probability mass under  $\mathbb{P}_x$  is (up to a multiplicative factor) a lower bound on the quantity  $A_1$  (see Eq equation C.3). Similarly, the region  $A_2$  is a subset of the intersection  $G \cap B_w$  with the band of height  $\leq 2$ , and has probability mass that is (up to a multiplicative factor) an upper bound on the quantity  $A_2$  (see Eq equation C.4).

Let  $F(y; x)$  denote  $\frac{1}{2} \mathbb{1}_{y \in \text{sign}(hw; xi)g}$ . We may write  $x_v$  as  $|x_w| \text{sign}(x_v)$  and let  $G \subset \mathbb{R}^2$  such that  $\text{sign}(x_v) \text{sign}(hw; xi) = 1$  iff  $x \in G$ . Then,  $\text{sign}(x_v) \text{sign}(hw; xi) = \mathbb{1}_{x \in 2G} - \mathbb{1}_{x \in 2G}$ . We get

$$\begin{aligned} & \| \text{proj}_h r_w L(w) \|_2 = \\ & = \mathbb{E} \left[ \phi^0(|x_w|) (\mathbb{1}_{x \in 2G} - \mathbb{1}_{x \in 2G}) F(y; x) |x_v| \right] \\ & = \mathbb{E} \left[ \phi^0(|x_w|) \mathbb{1}_{x \in 2G} F(y; x) |x_v| \right] - \mathbb{E} \left[ \phi^0(|x_w|) \mathbb{1}_{x \in 2G} F(y; x) |x_v| \right] \end{aligned}$$

Let  $A_1 = \mathbb{E} \left[ \phi^0(|x_w|) \mathbb{1}_{x \in 2G} F(y; x) |x_v| \right]$  and  $A_2 = \mathbb{E} \left[ \phi^0(|x_w|) \mathbb{1}_{x \in 2G} F(y; x) |x_v| \right]$ . (See Figure 3.) Note that  $\mathbb{E}_{y|x} [F(y; x)] = \frac{1}{2} \mathbb{1}_{x \in [1-2; 1]}$ , where  $1-2 > 0$ . Therefore, we have that  $A_1 = \frac{1}{2} \mathbb{E} \left[ \phi^0(|x_w|) \mathbb{1}_{x \in 2G} |x_v| \right]$  and  $A_2 = \mathbb{E} \left[ \phi^0(|x_w|) \mathbb{1}_{x \in 2G} |x_v| \right]$ .

Note that due to Proposition C.1,  $\phi^0(|x_w|) = c$  for some constant  $c$  and  $\phi^0(|x_w|) = 0$  whenever  $|x_w| > 2$ . Therefore, if  $U_2$  is the band  $B_w(\epsilon=2) = \{x : |x_w| \leq 2\}$  we have

$$A_2 \leq \frac{c}{2} \mathbb{E}[\mathbb{1}_{x \in 2G} \mathbb{1}_{x \in U_2} |x_v|] \tag{C.1}$$

Moreover, for each individual  $x$ , we have  $\phi^0(|x_w|) \mathbb{1}_{x \in 2G} |x_v| \geq 0$ , due to the properties of  $\phi^0$  (Proposition C.1). Hence, for any set  $U_1 \subset \mathbb{R}^d$  we have that

$$A_1 \geq \frac{1}{2} \mathbb{E} \left[ \phi^0(|x_w|) \mathbb{1}_{x \in 2G} \mathbb{1}_{x \in U_1} |x_v| \right]$$

Setting  $U_1 = B_w(\epsilon=6) = \{x : |x_w| \leq 6\}$ , by Proposition C.1, we get  $\mathbb{E} \left[ \phi^0(|x_w|) \mathbb{1}_{x \in 2G} \mathbb{1}_{x \in U_1} |x_v| \right] \geq \frac{1}{2} \mathbb{E}[\mathbb{1}_{x \in 2G} \mathbb{1}_{x \in U_1}]$ .

$$A_1 \geq \frac{1}{2} \mathbb{E}[\mathbb{1}_{x \in 2G} \mathbb{1}_{x \in U_1} |x_v|] \tag{C.2}$$

We now observe that by the definitions of  $U_1, U_2$ , for any constant  $R > 0$ , there exist some constants  $c^0, c^{00} > 0$  such that if  $\epsilon = \tan^{-1} c^0 R$  (the points in  $\mathbb{R}^2$  where  $\mathbb{C}$  intersects either  $U_1$  or

$U_2$  have projections on  $\mathcal{D}$  that are ( $\epsilon = \tan \theta$ ) we have that

$$\begin{aligned} & \mathbb{E}[f(x_{v,j}) \mid x_{v,j} \in \mathcal{D}] \leq \mathbb{E}[f(x_{v,j}) \mid x_{v,j} \in \mathcal{D}^c] + \epsilon \mathbb{E}[f(x_{v,j}) \mid x_{v,j} \in \mathcal{D}] \\ & \mathbb{E}[f(x_{v,j}) \mid x_{v,j} \in \mathcal{D}] \leq \mathbb{E}[f(x_{v,j}) \mid x_{v,j} \in \mathcal{D}^c] + \epsilon \mathbb{E}[f(x_{v,j}) \mid x_{v,j} \in \mathcal{D}] \end{aligned}$$

By equations equation C.1 and equation C.2, we get the following bounds whose graphical representations can be found in Figure 3.

$$A_1 \leq \frac{c^0 R (1 - \epsilon^2)}{\epsilon} \mathbb{E}[f(x_{v,j}) \mid x_{v,j} \in \mathcal{D}^c] + \epsilon \mathbb{E}[f(x_{v,j}) \mid x_{v,j} \in \mathcal{D}] \quad (C.3)$$

$$A_2 \leq \frac{c^0 c^{00}}{\tan \theta} \mathbb{E}[f(x_{v,j}) \mid x_{v,j} \in \mathcal{D}^c] + \epsilon \mathbb{E}[f(x_{v,j}) \mid x_{v,j} \in \mathcal{D}] \quad (C.4)$$

So far, we have used no distributional assumptions. Now, consider the corresponding expectations under the target margin  $\mathcal{D}$  (which we assumed to be strongly log-concave).

$$\begin{aligned} I_1 &= \mathbb{E}_{\mathcal{D}}[f(x_{v,j}) \mid x_{v,j} \in \mathcal{D}^c] + \epsilon \mathbb{E}[f(x_{v,j}) \mid x_{v,j} \in \mathcal{D}] \\ I_2 &= \mathbb{E}_{\mathcal{D}}[f(x_{v,j}) \mid x_{v,j} \in \mathcal{D}^c] + \epsilon \mathbb{E}[f(x_{v,j}) \mid x_{v,j} \in \mathcal{D}] \end{aligned}$$

Any strongly log-concave distribution enjoys the “well-behaved” properties defined by (DKTZ20a), and therefore, if  $R$  is picked to be small enough, then  $\mathcal{D}_1$  and  $\mathcal{D}_2$  are of order  $(\epsilon^2)$  (due to upper and lower bounds on the two dimensional marginal density over within constant radius balls – aka anti-anticoncentration and anticoncentration). Moreover, by Proposition A.2, we have  $\mathbb{E}[f(x_{v,j}) \mid x_{v,j} \in \mathcal{D}_1]$  and  $\mathbb{E}[f(x_{v,j}) \mid x_{v,j} \in \mathcal{D}_2]$  are both of order  $(\epsilon^2)$ . Hence we have that there exist constants  $c_1^0, c_2^0 > 0$  such that for the conditional expectations we have

$$\begin{aligned} \mathbb{E}_{\mathcal{D}}[f(x_{v,j}) \mid x_{v,j} \in \mathcal{D}^c] &\leq c_1^0 \\ \mathbb{E}_{\mathcal{D}}[f(x_{v,j}) \mid x_{v,j} \in \mathcal{D}] &\leq c_2^0 \end{aligned}$$

By assumption, Property equation 3.3 holds and, therefore, if  $c_1^0 = 2, c_2^0 = 2 =: c_2$ , we get that

$$\begin{aligned} \mathbb{E}_{\mathcal{D}_x}[f(x_{v,j}) \mid x_{v,j} \in \mathcal{D}^c] &\leq c_1^0 = 2 \\ \mathbb{E}_{\mathcal{D}_x}[f(x_{v,j}) \mid x_{v,j} \in \mathcal{D}] &\leq c_2^0 = 2 \end{aligned}$$

Moreover, by Property equation 3.2, we have that (under the true margin  $\mathcal{D}_1$  and  $\mathcal{D}_2$ ) are both  $(\epsilon^2)$ . Hence, in total, we get that for some constants  $\epsilon_1, \epsilon_2$ , we have

$$A_1 \leq \epsilon_1 (1 - \epsilon^2) \text{ and } A_2 \leq \epsilon_2 \frac{1}{\tan \theta}$$

Hence, if we pick  $\epsilon = ((1 - \epsilon^2) \tan \theta)$ , we get the desired result.

## C.2 PROOF OF PROPOSITION 4.4

For the following all the probabilities and expectations are over  $\mathcal{D}$ . First we observe that

$$\begin{aligned} \mathbb{P}[y \neq \text{sign}(hw; x_i)] &= \mathbb{P}[y \neq \text{sign}(hw; x_i) \mid y = \text{sign}(hw; x_i)] + \mathbb{P}[y \neq \text{sign}(hw; x_i)] \\ &= \mathbb{P}[\text{sign}(hw; x_i) \neq \text{sign}(hw; x_i)] + \text{opt} \end{aligned}$$

Then, we observe that by assumption  $\mathcal{D}$  satisfies Property equation 3.2, we have

$$\mathbb{P}[|hw; x_{ij}| > \epsilon] \leq C_3$$

and that

$$\mathbb{P}[\text{sign}(hw; x_i) \neq \text{sign}(hw; x_i) \mid |hw; x_{ij}| > \epsilon] \leq \mathbb{P}[|hw; x_{ij}| > \frac{\epsilon}{\tan \theta}]$$

where  $v$  is some vector perpendicular to  $\mathcal{D}$ . Using Markov's inequality, we get

$$\mathbb{P}[|hw; x_{ij}| > \frac{\epsilon}{\tan \theta}] \leq \frac{(\tan \theta)^k}{\epsilon^k} \mathbb{E}[|hw; x_{ij}|^k]$$

But, by assumption that  $D_{\mathcal{X}^Y}$  satisfies Property equation 3.1, there is some constant  $C_1 > 0$  such that  $\mathbb{E}[|\langle \mathbf{v}, \mathbf{x} \rangle|^k] \leq (C_1 k)^{k/2}$ . Thus

$$\begin{aligned} \mathbb{P}[\text{sign}(\langle \mathbf{w}, \mathbf{x} \rangle) \neq \text{sign}(\langle \mathbf{w}^*, \mathbf{x} \rangle)] &\leq \mathbb{P}[|\langle \mathbf{w}, \mathbf{x} \rangle| \leq \sigma] \\ &\quad + \mathbb{P}[\text{sign}(\langle \mathbf{w}, \mathbf{x} \rangle) \neq \text{sign}(\langle \mathbf{w}^*, \mathbf{x} \rangle) \cap |\langle \mathbf{w}, \mathbf{x} \rangle| > \sigma] \\ &\leq C_3 \sigma + \frac{(C_1 k)^{k/2} (\tan \theta)^k}{\sigma^k}. \end{aligned}$$

By picking  $\sigma$  appropriately in order to balance the two terms (note that this is a different  $\sigma$  than the one in Lemma 4.3), we get the desired result.

## D PROOFS FROM SECTION 5

### D.1 PROOF OF THEOREM 5.1

We will follow the same steps as for proving Theorem 4.1. Once more, we draw a sufficiently large sample so that our testers are ensured to accept with high probability when the true marginal is indeed the target marginal  $D^*$  and so that we have generalization, i.e. the guarantee that any approximate minimizer of the empirical error (error on the uniform empirical distribution over the sample drawn) is also an approximate minimizer of the true error. The algorithm we use is once more Algorithm 1, but this time we make multiple calls for different parameters  $\sigma$  (and we run  $T_1$  with higher  $k$ , as we will see shortly) and reject if any of these calls rejects. If we accept, we output the output of the execution of Algorithm 1 with the minimum empirical error.

The main difference between the Massart noise case and the agnostic case is that in the former we were able to pick  $\sigma$  arbitrarily small, while in the latter we face a more delicate tradeoff. To balance competing contributions to the gradient norm, we must ensure that  $\sigma$  is at least  $\Theta(\text{opt})$  while also ensuring that it is not too large. And since we do not know the value of  $\text{opt}$ , we will need to search over a space of possible values for  $\sigma$  that is only polynomially large in relevant parameters (similar to the approach of (DKTZ20b)). In our case, we may sparsify the space  $(0, 1]$  of possible values for  $\sigma$  up to accuracy  $\Theta((\frac{\epsilon}{\sqrt{k}})^{1+1/k})$  and form a list of  $\text{poly}(k/\epsilon)$  possible values for  $\sigma$ , one of which will satisfy  $c_1 \sigma - \Theta((\frac{\epsilon}{\sqrt{k}})^{1+1/k}) \leq \text{opt} \leq c_1 \sigma$ . hence, we perform the same (testing-learning) process for each of the possible values of  $\sigma$  and get a list of candidate vectors which is still of polynomial size.

The final step is, again, to use Proposition 4.4, after running tester  $T_1$  with parameter  $k$  (Proposition 3.2) and tester  $T_2$  with appropriate parameters for each of the candidate weight vectors. We get that our list contains a vector  $\mathbf{w}$  with

$$\mathbb{P}_{D_{\mathcal{X}^Y}} [y \neq \text{sign}(\langle \mathbf{w}, \mathbf{x} \rangle)] \leq \text{opt} + c \cdot k^{1/2} \cdot \theta^{1-1/(k+1)},$$

where  $\langle \mathbf{w}, \mathbf{w}^* \rangle \leq \theta := c_2 \sigma$  for  $\sigma$  such that  $c_1 \sigma - \Theta((\frac{\epsilon}{\sqrt{k}})^{1+1/k}) \leq \text{opt} \leq c_1 \sigma$ .

$$\mathbb{P}_{D_{\mathcal{X}^Y}} [y \neq \text{sign}(\langle \mathbf{w}, \mathbf{x} \rangle)] \leq \text{opt} + c\sqrt{k} \cdot \frac{c_2}{c_1} \text{opt} + \Theta \left( \frac{\epsilon}{\sqrt{k}} \right)^{1+\frac{1}{k}} \cdot 1^{-\frac{1}{k+1}} \leq O(\sqrt{k} \cdot \text{opt}^{1-\frac{1}{k+1}}) + \epsilon.$$

However, we do not know which of the weight vectors in our list is the one guaranteed to achieve small error. In order to discover this vector, we estimate the probability of error of each of the corresponding halfspaces (which can be done efficiently, due to Hoeffding's bound) and pick the one with the smallest error. This final step does not require any distributional assumptions and we do not need to perform any further tests.

In order to obtain our  $\tilde{O}(\text{opt})$  quasipolynomial time guarantee, observe first that we may assume without loss of generality that  $\text{opt} \geq 1/d^C$  for some  $C$ ; if instead  $\text{opt} = o(1/d^2)$ , say, then a sample of  $O(d)$  points will with high probability be noiseless, and so simple linear programming will recover a consistent halfspace that will generalize. Moreover, we may assume that  $\text{opt} \leq 1/10$ , since otherwise achieving  $O(\text{opt})$  is trivial (we may output an arbitrary halfspace). Let us adapt our algorithm so that we run tester  $T_1$  (see Proposition 3.2) multiple times for all  $k = 1, 2, \dots, \lceil \log^2 d \rceil$  (this only changes our time and sample complexity by a  $\text{polylog}(d)$  factor). Then Proposition 4.4

holds for some  $k^*$  such that  $k^* \in [\log(1/\text{opt}), 2 \log(1/\text{opt})]$ , since the interval has length at least 1 (and therefore it contains some integer) and  $2 \log(1/\text{opt}) \leq 2C \log d \leq \log^2 d$  (for large enough  $d$ ). Therefore, by picking the best candidate we get a guarantee of order

$$\begin{aligned} \sqrt{k^*} \cdot \text{opt}^{1-1/k} &= \sqrt{k^*} \cdot \text{opt}^{-1/k} \text{opt} \\ &= \sqrt{k^*} \cdot 2^{\frac{1}{k} \log \frac{1}{\text{opt}}} \cdot \text{opt} \\ &\leq \frac{\text{opt}}{2 \log(1/\text{opt})} \cdot 2 \cdot \text{opt} \quad (\text{since } \log(1/\text{opt}) \leq k^* \leq 2 \log(1/\text{opt})) \\ &= \Theta(\text{opt}). \end{aligned}$$

This concludes the proof of Theorem 5.1.

## D.2 PROOF OF LEMMA 5.2

In the agnostic case, the proof is analogous to the proof of Lemma 4.3. However, in this case, the difference is that the random variable  $F(y, \mathbf{x}) = 1 - 2 \cdot \{y \neq \text{sign}(\langle \mathbf{w}^*, \mathbf{x} \rangle)\}$  does not have conditional expectation on  $\mathbf{x}$  that is lower bounded by a constant. Instead, we need to consider an additional term  $A_3$  corresponding to the part  $2 \cdot \{y \neq \text{sign}(\langle \mathbf{w}^*, \mathbf{x} \rangle)\}$  and the term  $A_1$  will not be scaled by the factor  $(1 - 2\eta)$  as in Lemma 4.3. Hence, with similar arguments we have that

$$\|\nabla_{\mathbf{w}} \mathcal{L}_\sigma(\mathbf{w})\|_2 \geq A_1 - A_2 - A_3,$$

where  $A_1 \geq \tilde{c}_1$ ,  $A_2 \leq \tilde{c}_2 \cdot \frac{\sigma}{\tan \theta}$  and (using properties of  $\ell'_\sigma$  as in Lemma 4.3 and the Cauchy-Schwarz inequality)

$$\begin{aligned} A_3 &= 2 \mathbb{E}[\ell'_\sigma(|\mathbf{x}_{\mathbf{w}}|) \cdot \{ \mathbf{x} \in \mathcal{G} \} \cdot \{y \neq \text{sign}(\langle \mathbf{w}, \mathbf{x} \rangle)\} \cdot |\mathbf{x}_{\mathbf{v}}|] \leq \\ &\leq \frac{2c}{\sigma} \cdot \mathbb{E}[\{ \mathbf{x} \in \mathcal{U}_2 \} \cdot \{y \neq \text{sign}(\langle \mathbf{w}, \mathbf{x} \rangle)\} \cdot |\mathbf{x}_{\mathbf{v}}|] \leq \\ &\leq \frac{2c}{\sigma} \cdot \frac{\text{opt}}{\mathbb{E}[\{ \mathbf{x} \in \mathcal{U}_2 \} \cdot (\mathbf{x}_{\mathbf{v}})^2]} \cdot \frac{\text{opt}}{\mathbb{E}[\{y \neq \text{sign}(\langle \mathbf{w}, \mathbf{x} \rangle)\}]} = \\ &= \frac{2c\sqrt{\text{opt}}}{\sigma} \cdot \frac{\text{opt}}{\mathbb{E}[\langle \mathbf{v}, \mathbf{x} \rangle^2 \mid \mathbf{x} \in \mathcal{U}_2] \cdot \mathbb{P}[\mathbf{x} \in \mathcal{U}_2]}. \end{aligned}$$

Similarly to our approach in the proof of Lemma 4.3, we can use the assumed properties equation 3.2 and equation 3.4 to get that

$$A_3 \leq \tilde{c}_3 \frac{\sqrt{\text{opt}}}{\sqrt{\sigma}},$$

which gives that in order for the gradient loss to be small, we require  $\text{opt} \leq \Theta(\sigma)$ .

## D.3 PROOF OF THEOREM 5.3

Before presenting the proof of Theorem 5.3, we prove the following Proposition, which is, essentially, a stronger version of Proposition 4.4 for the specific case when the target marginal distribution  $D^*$  is the standard multivariate Gaussian distribution. Proposition D.1 is important to get an  $O(\text{opt})$  guarantee for the case where the target distribution is the standard Gaussian.

**Proposition D.1.** *Let  $D_{\mathcal{X}\mathcal{Y}}$  be a distribution over  $\mathbb{R}^d \times \{\pm 1\}$ ,  $\mathbf{w}^* \in \arg \min_{\mathbf{w} \in \mathbb{S}^{d-1}} \mathbb{P}_{D_{\mathcal{X}\mathcal{Y}}}[y \neq \text{sign}(\langle \mathbf{w}, \mathbf{x} \rangle)]$  and  $\mathbf{w} \in \mathbb{S}^{d-1}$ . Let  $\theta \geq \angle(\mathbf{w}, \mathbf{w}^*)$  and suppose that  $\theta \in [0, \pi/4]$ . Then, for a sufficiently large constant  $C$ , there is a tester that given  $\delta \in (0, 1)$ ,  $\theta$ ,  $\mathbf{w}$  and a set  $S$  of samples from  $D_{\mathcal{X}}$  with size at least  $\frac{d}{\theta} \log \frac{1}{\delta} \cdot C$ , runs in time  $\text{poly}(\frac{1}{\theta}, d, \log \frac{1}{\delta})$  and with probability  $1 - \delta$  satisfies the following specifications:*

- If the distribution  $D_{\mathcal{X}}$  is  $\mathcal{N}(0, I_d)$ , the tester accepts.
- If the tester accepts, then we have:

$$\Pr_{\mathbf{x} \sim S} [\text{sign}(\langle \mathbf{w}^*, \mathbf{x} \rangle) \neq \text{sign}(\langle \mathbf{w}, \mathbf{x} \rangle)] \leq O(\theta)$$

*Proof.* The testing algorithm does the following:

1. **Given:** Integer  $d$ , set  $S \subset \mathbb{R}^d$ ,  $\mathbf{w} \in \mathbb{S}^{d-1}$ ,  $\theta \in (0, \pi/4]$  and  $\delta \in (0, 1)$ .
2. Let  $\text{proj}_{\perp \mathbf{w}} : \mathbb{R}^d \rightarrow \mathbb{R}^{d-1}$  denote the operator that projects a vector  $\mathbf{x} \in \mathbb{R}^d$  to its projection into the  $(d-1)$ -dimensional subspace of  $\mathbb{R}^d$  that is orthogonal to  $\mathbf{w}$ .
3. For  $i \in \{0, \pm 1, \dots, \pm \frac{\sqrt{2 \log \frac{1}{\theta}}}{\theta}\}$ 
  - (a)  $S_i \leftarrow \{\mathbf{x} \in S : \langle \mathbf{w}, \mathbf{x} \rangle \in [i\theta, (i+1)\theta]\}$
  - (b) If  $\frac{|S_i|}{|S|} > 2\theta$ , then reject.
  - (c) If  $\frac{1}{|S_i|} \mathbb{P}_{\mathbf{x} \in S_i} (\text{proj}_{\perp \mathbf{w}}(\mathbf{x}))(\text{proj}_{\perp \mathbf{w}}(\mathbf{x}))^T - I_{(d-1)} \text{op} > 0.1$ , reject.
4. If  $\frac{1}{|S|} \mathbb{P}_{\mathbf{x} \in S} |\langle \mathbf{w}, \mathbf{x} \rangle| > \sqrt{2 \log \frac{1}{\theta}} > 5\theta$ , then reject.
5. If reached this step, accept.

If the tester accepts, then we have the following properties for some sufficiently large constant  $C' > 0$ . For the following, consider the vector  $\mathbf{v} \in \mathbb{R}^d$  to be the vector that is perpendicular to  $\mathbf{w}$ , lies within the plane defined by  $\mathbf{w}$  and  $\mathbf{w}^*$  and  $\langle \mathbf{v}, \mathbf{w}^* \rangle \leq 0$ .

1.  $\mathbb{P}_{\mathbf{x} \sim S} [|\langle \mathbf{w}, \mathbf{x} \rangle| \in [\theta i, \theta(i+1)]] \leq C'\theta$ , for any  $i \in \{0, \pm 1, \dots, \pm \frac{1}{\theta} \sqrt{2 \log \frac{1}{\theta}}\}$ .
2.  $\mathbb{P}_{\mathbf{x} \sim S_i} |\langle \mathbf{v}, \mathbf{x} \rangle| > \frac{\theta}{\tan \theta} \cdot i \leq C'/i^2$ , for any  $i \in \{0, \pm 1, \dots, \pm \frac{1}{\theta} \sqrt{2 \log \frac{1}{\theta}}\}$ .
3.  $\mathbb{P}_{\mathbf{x} \sim S} |\langle \mathbf{w}, \mathbf{x} \rangle| \geq \frac{1}{2 \log \frac{1}{\theta}} \leq C'\theta$ .

Then, for  $k = \frac{1}{\theta} \sqrt{2 \log \frac{1}{\theta}}$  and  $\text{Strip}_i = \{\mathbf{x} \in \mathbb{R}^d : \langle \mathbf{w}, \mathbf{x} \rangle \in [\theta i, \theta(i+1)]\}$ , we have that

$$\begin{aligned} & \Pr_{\mathbf{x} \sim S} [\text{sign}(\langle \mathbf{w}, \mathbf{x} \rangle) \neq \text{sign}(\langle \mathbf{w}^*, \mathbf{x} \rangle)] \leq \\ & \sum_{i=-k}^k \mathbb{P}_{\mathbf{x} \sim S} [\mathbf{x} \in \text{Strip}_i] \cdot \mathbb{P}_{\mathbf{x} \sim S} [|\langle \mathbf{v}, \mathbf{x} \rangle| > \frac{\theta}{\tan \theta} \cdot i \mid \mathbf{x} \in \text{Strip}_i] + \mathbb{P}_{\mathbf{x} \sim S} [|\langle \mathbf{w}, \mathbf{x} \rangle| \geq \frac{1}{2 \log \frac{1}{\theta}}] \leq \\ & \sum_{i=-k}^k \frac{|S_i|}{|S|} \cdot \mathbb{P}_{\mathbf{x} \sim S_i} [|\langle \mathbf{v}, \mathbf{x} \rangle| > \frac{\theta}{\tan \theta} \cdot i] + C'\theta \leq (C')^2 \theta \cdot \sum_{i \neq 0} \frac{1}{i^2} + C'\theta = O(\theta) \end{aligned}$$

Now, suppose the distribution  $D_{\mathcal{X}}$  is indeed the standard Gaussian  $\mathcal{N}(0, I_d)$ . We would like to show that our tester accepts with probability at least  $1 - \delta$ . Since  $D = \mathcal{N}(0, I_d)$ , we see that for  $\mathbf{x} \sim D$  we have that  $\mathbf{x} \cdot \mathbf{w}$  is distributed as  $\mathcal{N}(0, 1)$ . This implies that

- For all  $i \in \{0, \pm 1, \dots, \pm \frac{\sqrt{2 \log \frac{1}{\theta}}}{\theta}\}$  we have
  - $\Pr_{\mathbf{x} \sim \mathcal{N}(0, I_d)} [\langle \mathbf{w}, \mathbf{x} \rangle \in [i\theta, (i+1)\theta]] \leq \frac{1}{\sqrt{2\pi}} \theta$
  - $\Pr_{\mathbf{x} \sim \mathcal{N}(0, I_d)} [\langle \mathbf{w}, \mathbf{x} \rangle \in [i\theta, (i+1)\theta]] \geq \theta \cdot \min_{x \in [-\sqrt{2 \log \frac{1}{\theta}} - \theta, \sqrt{2 \log \frac{1}{\theta}} + \theta]} \frac{1}{\sqrt{2\pi}} e^{-\frac{x^2}{2}} \geq \frac{\theta^2}{10}$
- $\Pr_{\mathbf{x} \sim \mathcal{N}(0, I_d)} [\langle \mathbf{w}, \mathbf{x} \rangle \in [i\theta, (i+1)\theta]] \leq \frac{1}{\sqrt{2\pi}} \theta$
- $\Pr_{\mathbf{x} \sim \mathcal{N}(0, I_d)} |\langle \mathbf{w}, \mathbf{x} \rangle| > 2 \sqrt{2 \log \frac{1}{\theta}} = \int_{2\sqrt{2 \log \frac{1}{\theta}}}^{\infty} \frac{1}{\sqrt{2\pi}} e^{-\frac{x^2}{2}} dx \leq \theta \int_0^{\infty} \frac{1}{\sqrt{2\pi}} e^{-\frac{x^2}{2}} dx = \frac{\theta}{2}$

Therefore, via the standard Hoeffding bound, we see that for sufficiently large absolute constant  $C$  we have with probability at least  $1 - \frac{\delta}{4}$  over the choice of  $S$  that

- For all  $i \in \{0, \pm 1, \dots, \pm \frac{\sqrt{2 \log \frac{1}{\delta}}}{\theta}\}$  we have
  - $\Pr_{\mathbf{x} \sim S} [\langle \mathbf{w}, \mathbf{x} \rangle \in [i\theta, (i+1)\theta]] \leq \theta$
  - $\Pr_{\mathbf{x} \sim S} [\langle \mathbf{w}, \mathbf{x} \rangle \in [i\theta, (i+1)\theta]] \geq \frac{\theta^2}{20}$
- $\Pr_{\mathbf{x} \sim S} [\langle \mathbf{w}, \mathbf{x} \rangle > 2 \frac{\sqrt{2 \log \frac{1}{\delta}}}{\theta}] \leq \theta$
- $\Pr_{\mathbf{x} \sim S} [\langle \mathbf{w}, \mathbf{x} \rangle < -2 \frac{\sqrt{2 \log \frac{1}{\delta}}}{\theta}] \leq \theta$

Finally, we would like to show that conditioned on the above, the probability of rejection in step (3b) is small.

**Fact D.2.** *Given a set  $S \subset \mathbb{R}^{d-1}$  of i.i.d. samples from  $\mathcal{N}(0, I_d)$ , with probability at least  $1 - \text{poly} \frac{|S|}{d}$  we have*

$$\frac{1}{|S|} \sum_{\mathbf{x} \in S} \langle \mathbf{w}, \mathbf{x} \rangle \langle \mathbf{w}, \mathbf{x} \rangle^T - I_{(d-1)} \leq 0.1$$

Now, since each sample  $\mathbf{x}_i$  is drawn i.i.d. from  $\mathcal{N}(0, I_d)$ , we have that  $\langle \mathbf{w}, \mathbf{x}_i \rangle$  and  $\text{proj}_{\perp \mathbf{w}}(\mathbf{x}_i)$  are all independent from each other for all  $i$ . Since all the events we conditioned on depend on  $\{\langle \mathbf{w}, \mathbf{x}_i \rangle\}$  we see that  $\{\text{proj}_{\perp \mathbf{w}}(\mathbf{x}_i)\}$  are still distributed as i.i.d. samples from  $\mathcal{N}(0, I_{(d-1)})$ .

Recall that one of the events we have already conditioned on is that  $\Pr_{\mathbf{x} \sim S} [\langle \mathbf{w}, \mathbf{x} \rangle \in [i\theta, (i+1)\theta]] \geq \frac{\theta^2}{20}$  for all  $i \in \{0, \pm 1, \dots, \pm \frac{\sqrt{2 \log \frac{1}{\delta}}}{\theta}\}$ . This allows us to lower bound by  $\theta^2/20$  the ratio  $|S_i|/|S|$ . And since, as we described, for all these elements  $\mathbf{x}_i$  the vectors  $\text{proj}_{\perp \mathbf{w}}(\mathbf{x}_i)$  are distributed as i.i.d. samples from  $\mathcal{N}(0, I_{(d-1)})$ , we can use Fact D.2 to conclude that for sufficiently large absolute constant  $C$ , when  $|S| = \frac{d}{\theta} \log \frac{1}{\delta}^C$  we have with probability  $1 - \frac{\delta}{4}$  for all  $i \in \{0, \pm 1, \dots, \pm \frac{\sqrt{2 \log \frac{1}{\delta}}}{\theta}\}$  that

$$\frac{1}{|S_i|} \sum_{\mathbf{x} \in S_i} (\text{proj}_{\perp \mathbf{w}}(\mathbf{x})) (\text{proj}_{\perp \mathbf{w}}(\mathbf{x}))^T - I_{(d-1)} \leq 0.1$$

Overall, this allows us to conclude that with probability at least  $1 - \delta$  the tester accepts.  $\square$

We now present the proof of Theorem 5.3.

In the proof of Theorem 5.1, when the target distribution is the standard Gaussian in  $d$  dimensions, we may apply Proposition D.1 (and run the corresponding tester), instead of Proposition 4.4, in order to ensure that our list will contain a vector  $\mathbf{w}$  with

$$\begin{aligned} \mathbb{P}_{D_{\times Y}} [y \neq \text{sign}(\langle \mathbf{w}, \mathbf{x} \rangle)] &\leq \mathbb{P}_{D_{\times Y}} [y \neq \text{sign}(\langle \mathbf{w}^*, \mathbf{x} \rangle)] + \mathbb{P}_{D_{\times Y}} [\text{sign}(\langle \mathbf{w}^*, \mathbf{x} \rangle) \neq \text{sign}(\langle \mathbf{w}, \mathbf{x} \rangle)] \\ &\leq \text{opt} + O(\theta) \end{aligned}$$

where  $\langle \mathbf{w}, \mathbf{w}^* \rangle \leq \theta := c_2 \sigma$  and  $\sigma$  is such that  $c_1 \sigma - \Theta(\epsilon) \leq \text{opt} \leq c_1 \sigma$ , which gives the desired  $O(\text{opt}) + \epsilon$  bound. To get the value of  $\sigma$  with the desired property, we once again sparsified the space  $(0, 1]$  of possible values for  $\sigma$ , this time up to accuracy  $\Theta(\epsilon)$ .

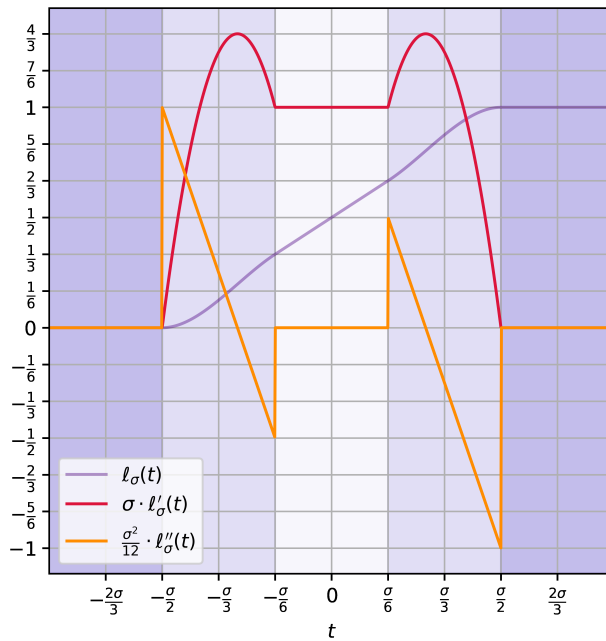


Figure 4: Figure illustrating the (normalized) first two derivatives of the function  $l_\sigma$  used to define the non convex surrogate loss  $\mathcal{L}_\sigma$ . The normalization is appropriate since  $l'_\sigma$  and  $l''_\sigma$  are homogeneous in  $1/\sigma$  and  $1/\sigma^2$  respectively. In particular, we see that  $l'_\sigma \leq \Theta(1/\sigma)$  and  $|l''_\sigma| \leq \Theta(1/\sigma^2)$  everywhere.