# Statistically Near-Optimal Hypothesis Selection

Olivier Bousquet\* Mark Braverman<sup>†</sup> Klim Efremenko<sup>‡</sup> Gillat Kol<sup>§</sup> Shay Moran<sup>¶</sup>

#### Abstract

Hypothesis Selection is a fundamental distribution learning problem where given a comparatorclass  $Q = \{q_1, \ldots, q_n\}$  of distributions, and a sampling access to an unknown target distribution p, the goal is to output a distribution q such that  $\mathsf{TV}(p,q)$  is close to  $\mathsf{opt}$ , where  $\mathsf{opt} = \min_i \{\mathsf{TV}(p,q_i)\}$  and  $\mathsf{TV}(\cdot,\cdot)$  denotes the total-variation distance. Despite the fact that this problem has been studied since the 19th century, its complexity in terms of basic resources, such as number of samples and approximation guarantees, remains unsettled (this is discussed, e.g., in the charming book by Devroye and Lugosi '00). This is in stark contrast with other (younger) learning settings, such as PAC learning, for which these complexities are well understood.

We derive an optimal 2-approximation learning strategy for the Hypothesis Selection problem, outputting q such that  $\mathsf{TV}(p,q) \leq 2 \cdot \mathsf{opt} + \varepsilon$ , with a (nearly) optimal sample complexity of  $\tilde{O}(\log n/\varepsilon^2)$ . This is the first algorithm that simultaneously achieves the best approximation factor and sample complexity: previously, Bousquet, Kane, and Moran (COLT '19) gave a learner achieving the optimal 2-approximation, but with an exponentially worse sample complexity of  $\tilde{O}(\sqrt{n}/\varepsilon^{2.5})$ , and Yatracos (Annals of Statistics '85) gave a learner with optimal sample complexity of  $O(\log n/\varepsilon^2)$  but with a sub-optimal approximation factor of 3.

We mention that many works in the *Density Estimation* (a.k.a., *Distribution Learning*) literature use Hypothesis Selection as a black box subroutine. Our result therefore implies an improvement on the approximation factors obtained by these works, while keeping their sample complexity intact. For example, our result improves the approximation factor of the algorithm of Ashtiani, Ben-David, Harvey, Liaw, and Mehrabian (*JACM '20*) for agnostic learning of mixtures of gaussians from 9 to 6, while maintaining its nearly-tight sample complexity.

<sup>\*</sup>Google Brain, Zürich. obousquet@google.com.

<sup>&</sup>lt;sup>†</sup>Department of Computer Science, Princeton University. mbraverm@princeton.edu.

<sup>&</sup>lt;sup>‡</sup>Department of Computer Science, Ben Gurion University. klimefrem@gmail.com.

Department of Computer Science, Princeton University. gillat.kol@gmail.com.

<sup>¶</sup>Department of Mathematics, Technion and Google Research. smoran@technion.ac.il.

## 1 Introduction

Hypothesis selection is a fundamental task in statistics, where a learner is getting a sample access to an unknown distribution p on some, possibly infinite, domain  $\mathcal{X}$ , and wishes to output a distribution q that is "close" to p. The problem was studied extensively over the last century and found many applications, most notably, in machine learning.

In this paper we study the hypothesis selection problem in the *agnostic* setting, where we assume a fixed finite<sup>1</sup> class  $\mathcal{Q}$  of reference distributions which is known to the learner, and which may or may not contain  $p^2$ . The goal of the learner is to output a distribution q that is at least as close to p as any of the distributions in  $\mathcal{Q}$  in *total variation* distance (denoted here  $\mathsf{TV}(\cdot, \cdot)$ ).

The statistical performance of a learner is measured using two parameters, denoted  $\alpha$  and  $m = m(n, \varepsilon, \delta)$ , where  $\alpha$  is the approximation factor of the algorithm and m is its sample complexity. Specifically, we say that a class of distributions  $\mathcal{Q} = \{q_1, \ldots, q_n\}$  is  $\alpha$ -learnable with sample complexity  $m(n, \varepsilon, \delta)$  if there is a (possibly randomized) learner such that for every  $\varepsilon, \delta > 0$  and every target distribution p, upon receiving  $m(n, \varepsilon, \delta)$  random samples from p, the learner outputs a distribution q satisfying  $\mathsf{TV}(p, q) \leq \alpha \cdot \min_{i \in [n]} \{\mathsf{TV}(p, q_i)\} + \varepsilon$  with probability at least  $1 - \delta$ . For the discussion below, we think of  $\delta$  as a small constant.

How good can a learner be? A-priori, it is not even clear that every class Q is learnable with finite sample complexity. Consider the following natural algorithm for hypothesis selection: estimate  $\mathsf{TV}(q_i,p)$  for every  $q_i \in Q$  and output the  $q_i$  that minimizes this quantity. While this algorithm clearly works (and even achieves an approximation factor of  $\alpha = 1$ ), estimating  $\mathsf{TV}(q_i,p)$  for any  $q_i$  requires  $\tilde{\Omega}(|\mathcal{X}|)$  samples from p (see, e.g., [JHW18]). Thus, if the domain  $\mathcal{X}$  is infinite (say  $\mathcal{X} = \mathbb{R}$ ), the sample complexity of this algorithm is not even finite. However, perhaps surprisingly, despite the impossibility of estimating the distance of p from even one of the distributions  $q_i$ , one can still find an approximate minimizer of the distances (even when  $\mathcal{X}$  is infinite!).

What are the smallest  $\alpha$  and m for which any given class of distributions  $\mathcal{Q}$  of size n is  $\alpha$ -learnable with sample complexity m? A seminal work by Yatracos [Yat85] (also see [DL96, DL97, DL01]) shows that any reference class Q of size n is 3-learnable with sample complexity  $O(\log n/\varepsilon^2)$ . For the case of n=2, Mahalanabis and Stefankovic [MS08] improve the approximation factor, constructing a 2-learner. This was extended by the recent work of Bousquet, Kane, and Moran [BKM19] to give a 2-approximation for any finite n, using a very different scheme. A matching lower bound of 2 on the approximation factor follows from the work of [CDSS14].

Although the work of [BKM19] obtains the optimal approximation factor for the agnostic hypothesis selection problem, the sample complexity of their scheme is  $\tilde{O}(\sqrt{n}/\varepsilon^{2.5})$ , which is exponential in the sample complexity of Yatracos's algorithm<sup>3</sup>. Deriving optimal learners with efficient sample complexity is left as the main open problem in their work. In this paper, we give a novel 2-learner with (near) optimal sample complexity, getting the best of both worlds.

<sup>&</sup>lt;sup>1</sup>See discussion of the infinite case at the end of this section.

<sup>&</sup>lt;sup>2</sup>The setting where p is assumed to be in Q is called the *realizable* setting.

<sup>&</sup>lt;sup>3</sup>We note that [BKM19] also provide poly(log| $\mathcal{X}$ |, log  $n, \varepsilon^{-1}$ ) sample complexity bounds, which can be better than their general  $\tilde{O}(\sqrt{n}/\varepsilon^{2.5})$  bound for finite domains  $\mathcal{X}$ .

**Density Estimation.** Hypothesis selection, and, in particular, Yatracos's algorithm, found applications beyond learning finite classes. Specifically, it is used as a basic subroutine in density estimation tasks where the goal is to learn an infinite class of distributions, in the realizable or agnostic setting<sup>4</sup>. A popular method, where the reference class  $\mathcal{Q}$  may be infinite, is the *cover method* (a.k.a. the skeleton method). In this method, one "covers" the class  $\mathcal{Q}$  by a finite  $\alpha$ -cover; that is, a subclass  $\mathcal{Q}' \subseteq \mathcal{Q}$  of distributions such that for every  $q \in \mathcal{Q}$  there exists  $q' \in \mathcal{Q}'$  with  $\mathsf{TV}(q,q') \leq \alpha$ . Often times it is the case that even if  $\mathcal{Q}$  is infinite, a finite  $\varepsilon$ -net  $\mathcal{Q}'$  exists, and Yatracos's agnostic learning algorithm can be applied on  $\mathcal{Q}'$  (see [DL01, Dia16] and references within for many such examples).

While the minimal possible size of such a cover Q' is often exponential in the natural parameters of the class  $Q^5$ , because Yatracos's algorithm has poly-logarithmic sample complexity, the obtained density estimation algorithm has a polynomial sample complexity. Since many density estimation results follow the cover method, or other related methods<sup>6</sup> that use Yatracos's algorithm as a black box, our algorithm can imply an improvement for all of these results. (We mention a couple of such examples below, in Section 1.4).

We note that in the realizable setting for density estimation, where the distribution p we wish to learn is in the infinite class Q of distributions we are considering (that is,  $\mathsf{opt} = 0$ ), one can typically get a better approximation factor by taking a finer cover (smaller  $\alpha$ ). By taking an  $\alpha$ -cover of Q, the above method results in a distribution q with  $\mathsf{TV}(p,q) \leq \alpha + 3\mathsf{opt} = \alpha$ . However, in the agnostic setting, even if we take a very small  $\alpha$ , the resulting  $\mathsf{TV}(p,q)$  may not be small as it is dominated by 3opt. By using the result of this paper in lieu of Yatracos's learning algorithm, this distance can be made 2opt.

#### 1.1 Our Results

We design a 2-learner for the agnostic hypothesis selection problem with sample complexity whose dependence on both n and  $\varepsilon$  is (near) optimal.

**Theorem 1.** Let Q be a finite class of distributions and let n = |Q|. Then, Q is 2-learnable with sample complexity<sup>7</sup>  $m(n, \varepsilon, \delta) = \tilde{O}\left((\log n \cdot \min(\log n, \log(1/\delta)) + \log(1/\delta))/\varepsilon^2\right)$ . In particular, for constant  $\delta > 0$ ,

$$m(n, \varepsilon, \delta) = \tilde{O}\left(\frac{\log n}{\varepsilon^2}\right).$$

Our learner in Theorem 1 is deterministic, and, as in the case for [BKM19], it only makes statistical queries. That is, our learner can be implemented in the restricted model where instead of getting random samples from p, the learner has access to an oracle that on a query  $(f, \varepsilon)$ 

<sup>&</sup>lt;sup>4</sup>In fact, learning infinite classes was a part of Yatracos's original motivation.

<sup>&</sup>lt;sup>5</sup>One easy example of an exponential cover is when  $\mathcal{Q}$  is the set of all convex combinations of k fixed distributions  $p_1, \ldots, p_k, i.e., \mathcal{Q} = \{\sum_{i \in [k]} \beta_i p_i : \sum_{i \in [k]} \beta_i = 1, \beta_i \geq 0\}$ . The set  $\mathcal{Q} = \{\sum_{i \in [k]} \frac{r_i}{\ell} \cdot p_i : r_i \in \mathbb{N} \cup \{0\}, \ \ell = \lceil \frac{k}{\alpha} \rceil, \ \sum_{i \in [k]} \frac{r_i}{\ell} = 1\}$  is a cover of  $\mathcal{Q}$  of exponential size (in k). Sub-exponential covers are not possible in this case. See Chapter 7.4 in [DL01] for this example, and the rest of Chapter 7 for more such examples.

<sup>&</sup>lt;sup>6</sup>Another such method is the recent sample compression method by [ABDH<sup>+</sup>20], used to obtain improved density algorithms for the mixtures of Gaussians problem.

<sup>&</sup>lt;sup>7</sup>We use the standard notation that  $f(n) = \tilde{O}(h(n_1, \dots, n_t))$  if there exists  $k \in \mathbb{N}$  such that  $f(n_1, \dots, n_t) = O(h(n_1, \dots, n_t) \log^k(h(n_1, \dots, n_t)))$ .

answers by a value in  $\mathbb{E}_{x\sim p}[f(x)] \pm \varepsilon$  (or, equivalently, on a query  $(F,\varepsilon)$ , where F is a set, answers by  $p(F) \pm \varepsilon$ ). Furthermore, our algorithm consists of only  $\tilde{O}(\log n/\varepsilon^2)$  such rounds of queries, whereas the algorithm [BKM19] consists of  $O(n/\varepsilon)$  such rounds.

### 1.2 Our Technique

#### 1.2.1 The Cutting-With-Margin Game

To prove Theorem 1, we reduce the hypothesis selection problem to solving a geometric game we call the "cutting-with-margin" game. This game is between a player and an adversary and it is played over a convex body  $\mathcal{H} \subseteq \Delta_n$  known to both parties, where  $\Delta_n$  denotes the simplex of n-dimensional probability vectors<sup>8</sup>. In every round of the game, the player selects a point  $h \in \mathcal{H}$  and adversary updates the set  $\mathcal{H}$  to a new convex set by "cutting out" a part of  $\mathcal{H}$  that contains the  $\ell_1$  ball of radius  $\varepsilon$  around h. The game ends when the set  $\mathcal{H}$  is empty.

We first show that any strategy for the player which ensures that the game ends in at most r rounds implies a 2-learner for the hypothesis selection problem with sample complexity  $\tilde{O}(r \log n/\varepsilon^2)$  (this is because the implementation of each round requires n statistical queries that should be approximated to within  $O(\varepsilon)$ ). We then give an information-theoretic argument showing that the game is solvable in  $r = \tilde{O}(\log(n)/\varepsilon^2)$  rounds, implying a hypothesis selection algorithm with  $\tilde{O}(\log^2(n)/\varepsilon^4)$  samples. Our player's strategy views each point  $h \in \mathcal{H} \subseteq \Delta_n$  as a distribution and takes the point  $h \in \mathcal{H}$  that maximizes the entropy function.

Even though the cutting-with-margin game serves as a technical tool in this work, this simple game may also be of independent interest, and it is natural to study it for different norms (other than the  $\ell_1$  norm considered in this paper). In a sense, this game is a dual perspective on the geometric approach taken by [BKM19] (see Section 2). Nevertheless, it is the move to this dual perspective that allowed us to use the above maximum-entropy-based strategy. While entropy-based strategies are widely used in online optimization (see Section 1.4), we find the fact that such a strategy is helpful for making progress in this abstract statistical problem of hypothesis selection, to be curious. We hope that this connection will inspire more collaboration between the optimization and the statistical learning communities.

#### 1.2.2 Achieving Optimal Sample Complexity

Our solution for the cutting-with-margin game yields a hypothesis selection algorithm with sample complexity polynomial in  $\log n/\varepsilon$ , but still sub-optimal. While reducing the sample complexity of this algorithm and achieving a near optimal complexity of  $\tilde{O}(\log n/\varepsilon^2)$  requires quite a bit of effort (in fact, it is the main technical contribution of this paper), we believe that it makes our algorithm more applicable (in the sense that it can replace Yatracos's algorithm, without compromising the sample complexity).

To this end, at a very high level, we consider a "dynamic" cutting-with-margin game that allows the cutting of  $\ell_1$  balls of different diameters, and we give a "win-win"-style strategy, where in rounds where we use more samples the diameter of the ball we cut is larger (see Section 2.4). Thus, the player either makes a lot of progress towards the goal or uses few samples.

<sup>&</sup>lt;sup>8</sup>*I.e.*,  $\Delta_n := \{ h \in \mathbb{R}^n : \sum_{i \in [n]} h_i = 1, \ (\forall i) : h_i \ge 0 \}.$ 

A detailed overview of our techniques can be found in Section 2.

Adaptive data analysis. As explained in Section 2, the ("primal") geometric approach of [BKM19] results in a hypothesis selection algorithm that makes  $O(n^2/\varepsilon)$  statistical queries, where each should be approximated to within  $O(\varepsilon)$ . Had all these queries been submitted together, the standard combination of Chernoff and union bound would imply a logarithmic sample complexity. However, their algorithm submits these queries adaptively, in  $O(n/\varepsilon)$  rounds, where in each round n queries are submitted. Thus, naively, each of the rounds will require  $\tilde{O}(\log n/\varepsilon^2)$  fresh samples for the total sample complexity of  $\tilde{O}(n/\varepsilon^3)$ . Their improved stated sample complexity of  $\tilde{O}(\sqrt{n}/\varepsilon^{2.5})$  is made possible by importing clever tools from Adaptive Data Analysis.

Given the above, a natural question is whether similar "off-the-shelf" Adaptive Data Analysis tools can be used to convert the hypothesis selection algorithm obtained in Section 1.2.1 from our solution of the cutting-with-margin game, to a sample optimal one. (Recall that this protocol consists of  $\tilde{O}(\log n/\varepsilon^2)$  rounds and makes n statistical queries in each round). Unfortunately, we were unable to apply these tools to get a significant quantitative improvements, as these tools are mostly geared toward cases where there are many rounds of adaptivity, while in our algorithm, the number of rounds  $\tilde{O}(\log n/\varepsilon^2)$  is much smaller than the number of queries n made in every round (see, e.g., [DFH<sup>+</sup>15]). Instead, as described above, we use a more direct solution and tune the number of samples we use for each query adaptively, by monitoring (and verifying) the progress of the algorithm.

It will be interesting to explore whether our technique can be extended to more general protocols in adaptive data analysis.

### 1.3 Additional Discussion of The Model

In this work, we give an *improper* algorithm for the *finite agnostic* hypothesis selection problem under the *total variation distance*. We next explain the modeling choices we have made:

The finite agnostic setting. We consider the finite agnostic setting; clearly, an algorithm in this setting applies in the realizable setting as well. In addition, as discussed above, hypothesis selection in the finite agnostic setting is often used as a building block in the infinite (agnostic and realizable) settings (*i.e.*, in density estimation).

**Total variation distance.** The total variation distance is used by numerous prior works in the field, and is a natural choice for our study for several reasons: firstly, solving the hypothesis selection problem for the total variation distance (which corresponds to the  $\ell_1$  norm) implies solving the corresponding problem for any  $\ell_p$  norm, for  $p \in [1, \infty]$ , as  $||x - y||_p \le ||x - y||_1$ . Another reason is that for many other metrics, the sample complexity of a hypothesis selection problem can depend on structural properties of the reference class  $\mathcal{Q}$ , which is undesirable for formulating problem-independent theorems like Theorem 1. For a more elaborate discussion of the advantages in working with total variation, see Chapter 6.5 in [DL01], and Section 3.1 in [ABDH<sup>+</sup>20].

We believe that our technique can be extended to derive hypothesis selection algorithms for

other distance measures that satisfy (at least some approximate) version of the triangle inequality (e.g., Hellinger distance and other metric spaces).

**Proper vs. improper.** A basic classification of machine learning problems distinguishes between *proper* and *improper* learning. In the proper case the algorithm always outputs a distribution  $q \in \mathcal{Q}$ , whereas in the improper case it may output an arbitrary distribution. Improperness has been shown to be beneficial in many settings (see, e.g., [SF12, DS14]), including the agnostic hypothesis selection setting: while Yatracos's 3-approximation algorithm is proper, [BKM19] prove that the factor 3 cannot be improved by any proper algorithm (with any sample complexity)<sup>10</sup>. For this reason, their and our 2-approximation algorithms are inherently improper. For many applications (e.g., applications to density estimation discussed above), improper hypothesis selection algorithms suffice.

Computational complexity. Although our approach is algorithmic, our focus is not on computational efficiency. While the sample complexity of our algorithm is only logarithmic in the number of distributions n (and is independent of the domain size  $|\mathcal{X}|$ ), in the general case, its running time scales polynomially with both n and  $|\mathcal{X}|$ , as is the case for other sample-efficient hypothesis selection algorithms. Clearly, the dependence on n cannot be sub-linear (each  $q_i$  needs to be accessed, unless some structure on  $\mathcal{Q}$  is assumed). As for the dependence on  $|\mathcal{X}|$ , our algorithm assumes oracle access to operations on  $\mathcal{X}$ , such as checking membership in sets of the form  $F = \{x \in \mathcal{X} : q_1(x) > q_2(x)\}^{11}$ , and several other (somewhat involved) operations n that can only be implemented efficiently for restricted classes  $\mathcal{Q}$ . We mention that the situation is similar for many density estimation problems: the existence of polynomial time algorithms is unknown even for specific natural classes, such as mixtures of gaussians (see [ABDH+20] for further discussion).

While efficient algorithms (e.g., with poly  $\log(|\mathcal{X}|)$  running-time) for all classes  $\mathcal{Q}$  are unlikely in the simple and abstract learning setting considered by this work, this setting is particularly suited to capture basic information-theoretic resources, such as sample-complexity and approximation guarantees, which are not affected by the computational model. As discussed above, the complexity of these resources is still poorly understood, even for very basic problems.

#### 1.4 Additional Related Work

In this work we give a novel approximation algorithm for hypothesis selection of any (finite) class Q, following the classical work of [Yat85, DL96, DL97, DL01] and the recent work of [BKM19], discussed above. Over the last decade or so, hypothesis selection received quite a bit of attention by different theoretical communities and many aspects of this problem were studied, including computational efficiency, robustness, weaker access to hypotheses, privacy and more (see,

<sup>&</sup>lt;sup>9</sup>See Section 2.1 for our usage of the triangle inequality.

 $<sup>^{10}</sup>$ We mention that for the case n=2, a proper 2-approximation algorithm for the agnostic hypothesis selection problem was given by [MS08].

<sup>&</sup>lt;sup>11</sup>These are the, so called, "Yatracos sets" and Yatracos's algorithm also assumes membership oracle to them.

<sup>&</sup>lt;sup>12</sup>In the language of the overview presented in Section 2, these operations include finding a distribution q such that  $v(q) \leq v$ , and solving the optimization problem corresponding to finding the discriminating sets  $F_i$ .

e.g., [MS08, DDS15, DK14, SOAJ14, AJOS14, CDSS14, DKK<sup>+</sup>19, BKSW21, AFJ<sup>+</sup>18, BKSW21, GKK<sup>+</sup>20]).

Hypothesis selection can also be viewed as a special case of density estimation (also known as distribution learning), where one wishes to learn a (typically infinite) class of densities from samples. In fact, as mentioned above, many density estimation algorithms use hypothesis selection algorithms as fundamental subroutines. Density estimation is a very basic unsupervised learning problem studied since the late nineteenth century, starting with the pioneering work of Pearson [Pea95]. Since, it was systematically studied for many natural classes, such as mixtures of gaussians (e.g., [KMV12, DKS17, DKS18, KSS18, ABM18, ABDH+20]), histograms (e.g., [Pea95, LN96, DL04, CDSS14, DLS18]), and more. For a fairly recent survey see [Dia16].

Our result yields improved approximation guarantees in many of these works. For example, plugging it in [ABDH<sup>+</sup>20], instead of Yatracos's algorithm which is used as a black box, improves the approximation factor from 3 to 2 for learning gaussians, and from 9 to 6 for learning mixtures of gaussians, while keeping the sample complexity near-optimal.

Optimization and online learning. A key component in our derivation is the cutting-with-margin game. This game is reminiscent of dynamical processes which are studied in optimization and online learning. In particular, our solution to this game is based on a greedy approach of maximizing the entropy and a potential-based analysis which brings to mind standard KL-divergence-based analyses of mirror-decent and multiplicative-weights update (see, e.g., [AW01, AHK12, Bub15]). Moreover, the cutting-with-margin game naturally generalizes to arbitrary norms  $\|\cdot\|$  by replacing the  $\ell_1$  norm with  $\|\cdot\|$  and the simplex  $\Delta_n$  by the unit ball with respect to  $\|\cdot\|$ . One can extend our upper bound to arbitrary norms, by replacing the KL-divergence with an appropriate  $Bregman\ divergence^{13}$ , as is the case for some optimization problems.

These technical interrelations suggest the possibility of a deeper connection between the cutting-with-margin game and online optimization. Ideally, one could hope to find a formal reduction by phrasing our game as a convex regret minimization problem. We remark, however, that, unlike regret minimization problems, our game is not defined via a local regret function, but rather defined using a very global cost function. We leave this further exploration of the relations between our game to the regret minimization framework for future work.

The ellipsoid method. Another known algorithm that is of a particular syntactic similarity to our cutting-with-margin game is the well-known ellipsoid method for solving linear programs: in both settings a player maintains a convex set in  $\mathbb{R}^n$  (in our game it is, without loss of generality, a polytope, and when running the ellipsoid method it is an ellipsoid), and in each step it selects a point within that set. If the selected point is not a "solution", the player receives a separating hyperplane from an adversary or a hyperplane oracle, which separates the selected point from the target set of solutions. Then, the player moves to a "smaller" convex body that lies, in its entirety, on one side of the hyperplane.

<sup>&</sup>lt;sup>13</sup>Using the Bregman divergence, we have some preliminary results regarding the round complexity of our cutting-with-margin game in other norms. These include a nearly tight bounds for the  $\ell_p$  norm, when  $p \in (1,2] \cup \{\infty\}$ : if  $p \in (1,2)$  then the player can solve the corresponding game in  $r = O_p(1/\varepsilon^2)$  rounds, and if  $p = \infty$  a then the round complexity of the game is  $\Theta(n \log(1/\varepsilon))$ .

We note that a crucial difference between the two is that when running the ellipsoid method, the ellipsoids are getting rapidly smaller in terms of *volume* (and, for example, the next ellipsoids need not be contained in the former one), and it is this decrease in volume that allows for a fast convergence. In contrast, as will be discussed in Section 2.3, shrinking the volume of our convex body between rounds of the cutting-with-margin game does not suffice for convergence (and therefore, "centroid-based" methods do not apply).

### 2 Proof Overview

In this section we overview the proofs and highlight some of the more technical arguments. We defer the full proof to the Appendix.

Let  $Q = \{q_1, \ldots, q_n\}$  be a (known) finite reference class of distributions and let p denote the target distribution to which we have sample access. Denote  $i^* = \arg\min_i \{\mathsf{TV}(p, q_i)\}$ . Our goal is to use as few samples as possible from p in order to find q such that  $\mathsf{TV}(p, q) \leq 2 \cdot \mathsf{TV}(p, q_{i^*}) + \varepsilon$ .

## 2.1 A Geometric Approach to Hypothesis Selection

Our starting point is the 2-approximation algorithm of [BKM19]. In this subsection we describe our interpretation of their technique (some of the claims we make here are implicit in their paper).

The basic observation of [BKM19] is that it suffices to find a distribution q which is (almost) at least as close to each of the  $q_i$ 's as p,

$$(\forall i): \mathsf{TV}(q, q_i) \le \mathsf{TV}(p, q_i) + \varepsilon. \tag{1}$$

Finding such a q suffices, as by the triangle inequality,  $\mathsf{TV}(q,p) \leq \mathsf{TV}(q,q_i) + \mathsf{TV}(q_i,p) \leq 2\mathsf{TV}(q_i,p) + \varepsilon$  for every i, and, in particular, for  $i^*$ .

This suggests the following definitions: for a distribution q, let  $v(q) \in [0,1]^n$  denote the vector of all distances  $v(q) = (\mathsf{TV}(q,q_i))_{i=1}^n$ ; a vector  $v \in [0,1]^n$  is feasible if  $v \geq v(q)$  for some distribution q (when we write  $u \geq w$  for  $u, w \in [0,1]^n$  we mean  $(\forall i) : u_i \geq w_i$ ). With this notation, our goal is to find v such that

- (i)  $v \leq v(p) + \varepsilon \cdot 1_n$ , where  $1_n$  is the all-one vector, and
- (ii) v is feasible.

Once such a vector v is obtained, one can find a distribution q satisfying  $v(q) \leq v$ , and consequently a 2-approximation for the target distribution p.

Let  $\mathcal{P} \subseteq [0,1]^n$  denote the set of all feasible vectors v and note that it is convex and upwardclosed. The approach of [BKM19] for finding a desired v proceeds in rounds, where in round k we find a vector  $u_k$  that is closer to the feasible set, while maintaining the invariant that  $u_k \leq v(p)$ :

- 1. Let  $u_0 = \vec{0} \in [0,1]^n$  be the all-zero vector. Note that  $u_0 \leq v(p)$ , so  $u_0$  satisfies the above Item (i), but not Item (ii) (except in trivial cases).
- 2. For k = 0, 1, ...
  - (a) If  $u_k + \varepsilon \cdot 1_n$  is feasible (that is, if  $d_{\infty}(u_k, \mathcal{P}) \leq \varepsilon$ , where  $d_{\infty}(\cdot, \cdot)$  denotes  $\ell_{\infty}$  distance), then output a q such that  $v(q) \leq u_k + \varepsilon \cdot 1_n$  ( $\leq v(p) + \varepsilon \cdot 1_n$ ).

(b) Else, use samples from p to derive  $u_{k+1}$  such that  $u_k \leq u_{k+1} \leq v(p)$ , and  $u_{k+1}$  is "closer" (in some measure, see below) to  $\mathcal{P}$ .

Selecting the new point  $u_{k+1}$ . The crux of this approach is the update step in which  $u_{k+1}$  is computed given  $u_k$ . Since  $d_{\infty}(u_k, \mathcal{P}) > \varepsilon$ , there exists a  $u_{k+1}$  such that  $u_k \leq u_{k+1} \leq v(p)$  and  $d_1(u_{k+1}, u_k) \geq \frac{\varepsilon}{2}$  (for instance, since there exists a coordinate  $i \in [n]$  such that  $u_k + \frac{\varepsilon}{2} \cdot e_i < v(p)$ , where  $e_i$  is the  $i^{\text{th}}$  unit vector). [BKM19] show how to find such a  $u_{k+1}$  with few queries (discussed next), and they use this  $u_{k+1}$  as their next point. However, since  $||1_n||_1 = n$ , their strategy may require  $\Omega(\frac{n}{\varepsilon})$  rounds.

#### 2.1.1 Implementing the Strategy

Violated tests. We next explain how [BKM19] find the coordinate i of  $u_k$  that they wish to update. To this end, observe that whenever  $u_k + \varepsilon \cdot 1_n$  is not feasible there is a hyperplane separating the point  $u_k + \varepsilon \cdot 1_n$  from the set  $\mathcal{P}$  of feasible vectors, witnessing the fact that  $d_{\infty}(u, \mathcal{P}) > \varepsilon$ . We call a normal  $h \in \Delta_n$  to such a hyperplane a "violated test" (here  $\Delta_n$  denotes the simplex of all probability vectors in  $\mathbb{R}^n$ ). For  $u \in [0,1]^n$  and d > 0, we denote the set of all violated tests witnessing the fact that  $u + d \cdot 1_n$  is not feasible by

$$\mathcal{H}_d(u) = \Big\{ h \in \Delta_n : \ h \cdot u + d < \min_{v \in \mathcal{P}} h \cdot v \Big\}.$$

From a test h to an updated point  $u_{k+1}$ . We next informally state a central lemma proved by [BKM19], showing how to convert any violated test h to a new point  $u_{k+1}$  (for a precise statement, see Lemma 12 in [BKM19] or Lemma 7 in this paper).

**Lemma 2.** Using n statistical queries (queries of the form p(F) for some set F), any  $h \in \mathcal{H}_{\varepsilon}(u_k)$  can be converted to a point  $u_{k+1}$  satisfying:

- 1.  $u_k \le u_{k+1} \le v(p)$ .
- 2.  $u_{k+1}$  passes the test induced by  $h: h \notin \mathcal{H}_{\frac{\varepsilon}{2}}(u_{k+1})$ . This also implies that  $h \cdot (u_{k+1} u_k) > \frac{\varepsilon}{2}$  (as  $h \in \mathcal{H}_{\varepsilon}(u_k)$  implies  $h \cdot u_k + \varepsilon < \min_{v \in \mathcal{P}} h \cdot v$  and  $h \notin \mathcal{H}_{\frac{\varepsilon}{2}}(u_{k+1})$  implies  $h \cdot u_{k+1} + \frac{\varepsilon}{2} \ge \min_{v \in \mathcal{P}} h \cdot v$ ).

Observe that the  $u_{k+1}$  constructed by this lemma (for any h) satisfies  $d_1(u_{k+1}, u_k) \ge \frac{\varepsilon}{2}$  (due to Item 2, recall that  $h \in \Delta_n$ ), and therefore it can be used to implement the strategy of [BKM19].

**Proving the lemma.** While the proof of Lemma 2 is pretty short, it is tricky. For completeness, we will next give some intuition for it by showing how to construct  $u_{k+1}$  for a specific (easy to handle) h.

Assume that  $u_k + \varepsilon \cdot 1_n$  is not feasible and that  $h = (\frac{1}{2}, \frac{1}{2}, 0, \dots, 0) \in \mathcal{H}_{\varepsilon}(u_k)$ . Denote  $F = F(q_1, q_2) = \{x : q_1(x) \ge q_2(x)\}$ . (Observe that this is the so-called Yatracos set which is used in Yatracos's 3-approximation algorithm and satisfies  $\mathsf{TV}(q_1, q_2) = q_1(F) - q_2(F)$ ). Use samples from p to get an estimate  $\hat{p}(F)$  of p(F) up to an  $\frac{\varepsilon}{4}$  additive term. Set  $z_i = |\hat{p}(F) - q_i(F)| - \frac{\varepsilon}{2}$  for i = 1, 2 and  $z_i = 0$  for  $i \ge 3$ . Obtain  $u_{k+1}$  from  $u_k$  by setting  $(u_{k+1})_i = \max\{(u_k)_i, z_i\}$ .

The resulting  $u_{k+1}$  satisfies Item 1, as since  $|p(F) - q_i(F)| \le \mathsf{TV}(p, q_i) = (v(p))_i$  it follows that  $z_i \le (v(p))_i$ . It also satisfies Item 2, as

$$h \cdot u_{k+1} + \frac{\varepsilon}{2} = \frac{1}{2}((u_{k+1})_1 + (u_{k+1})_2) + \frac{\varepsilon}{2} \ge \frac{1}{2}(z_1 + z_2) + \frac{\varepsilon}{2}$$

$$\ge \frac{1}{2}(|\hat{p}(F) - q_1(F)| + |\hat{p}(F) - q_2(F)|) \ge \frac{1}{2}|q_1(F) - q_2(F)| = \frac{1}{2}\mathsf{TV}(q_1, q_2) = \min_{v \in \mathcal{P}} h \cdot v,$$
(2)

where the last equality is because for every  $v = v(q) \in \mathcal{P}$  it holds that  $h \cdot v = \frac{1}{2}(v_1 + v_2) = \frac{1}{2}(\mathsf{TV}(q, q_1) + \mathsf{TV}(q, q_2)) \ge \frac{1}{2}\mathsf{TV}(q_1, q_2)$  and for  $v = v(q_1) \in \mathcal{P}$  it holds that  $h \cdot v = \frac{1}{2}\mathsf{TV}(q_1, q_2)$ .

Query/sample complexity. For a general h, the proof of the lemma is more involved and crucially relays on the Minmax theorem. The point  $u_{k+1}$  is computed as  $(u_{k+1})_i = \max\{(u_k)_i, z_i\}$ , where for every  $i \in [n]$ ,  $z_i$  is of the form  $z_i = |\hat{p}(F_i) - q_i(F_i)| - \frac{\varepsilon}{2}$ , for some set  $F_i$  and where  $\hat{p}(F_i)$  is an approximation of  $p(F_i)$  to within an additive error of  $c \cdot \varepsilon$  for some constant c < 1.

Computing  $u_{k+1}$  requires n statistical queries (the values of  $p(F_i)$  for all i's), where each needs to be approximated to within an additive error of  $c \cdot \varepsilon$ . While approximating each query separately requires  $\Theta(1/\varepsilon^2)$  samples, by a standard combination of Chernoff and union bound, all n queries can be approximated using  $O(\log n/\varepsilon^2)$  samples.

### 2.2 The Cutting-With-Margin Game: A Dual Perspective

Recall that we wish to find a rule for updating  $u_k$  to a  $u_{k+1}$  satisfying  $u_k < u_{k+1} < v(p)$  that will allow us to reach a feasible point after the minimum number of steps. We wish to define a measure of progress to help us choose our next  $u_{k+1}$ . As discussed above, [BKM19] use the  $\ell_1$  norm as their measure of progress, but this results in a slow convergence to a feasible point.

To find a better progress measure, we revisit Lemma 2, specifically Item 2 that shows that by updating  $u_k$  using the test  $h \in \mathcal{H}_{\varepsilon}(u_k)$ , it is not only that  $h \notin \mathcal{H}_{\varepsilon}(u_{k+1})$ , but also  $h \notin \mathcal{H}_{\frac{\varepsilon}{2}}(u_{k+1})$ . We interpret this as implying that the set of violated tests can shrink substantially between rounds. This suggests a new approach: instead of measuring progress by comparing the locations of  $u_k$  and  $u_{k+1}$ , we can take a "dual" view and compare the sizes of the sets  $\mathcal{H}_{\varepsilon}(u_k)$  and  $\mathcal{H}_{\varepsilon}(u_{k+1})$  of violated tests that we still need to rule out (recall that if this set is empty, we have found a feasible point). We note that this "dual" view is lossy (and is not a dual in the standard sense) as the mapping  $u_k \to \mathcal{H}_{\varepsilon}(u_k)$  may not be one-to-one.

The cutting-with-margin game. Consider a sequence  $\vec{0} = u_0 \leq u_1 \leq \ldots \leq u_m$  in which the point  $u_{k+1}$  was produced from  $u_k$  by selecting some  $h_k \in \mathcal{H}_{\varepsilon}(u_k)$  and applying Lemma 2, and where  $u_m$  is feasible. Denote  $\mathcal{H}_k = \mathcal{H}_{\varepsilon}(u_k)$ . It can be shown that  $\mathcal{H}_k$  is convex for every k, and that  $\mathcal{H}_0 \supset \mathcal{H}_1 \supset \mathcal{H}_2 \supset \ldots \supset \mathcal{H}_m = \emptyset$  ( $\mathcal{H}_m = \emptyset$  as  $u_m$  is feasible). Furthermore, we are able to prove that  $\mathcal{H}_{k+1}$  is disjoint from an  $\ell_1$  ball of radius  $\Omega(\varepsilon)$  around  $h_k$  (see Lemma 9). Intuitively, this is because  $h_k \notin \mathcal{H}_{\frac{\varepsilon}{2}}(u_{k+1})$  (Lemma 2, Item 2) implies that the generated  $u_{k+1}$  not only passes the test induced by  $h_k$ , but also passes all "similar" tests.

The above discussion gives rise to the *cutting-with-margin* game discussed in the introduction (see Section 1.2.1). Recall that this is a game between a player and an adversary, and it is played over a convex body  $\mathcal{H} \subseteq \Delta_n$  known to both the player and the adversary. Let  $\mathcal{H}_0 = \mathcal{H}$ ; in

every round k = 0, 1, ... of the game, the player selects a point  $h_k \in \mathcal{H}_k$  and the adversary picks  $\mathcal{H}_{k+1} \subseteq \mathcal{H}_k$  to be any convex set which is disjoint from the  $\ell_1$  ball of radius  $\varepsilon$  around  $h_k$ . The game ends when the set  $\mathcal{H}_k$  is empty. See illustration in Figure 1. Of course, the task is now to find a strategy that solves this game with minimum number of rounds. Note that, in the language of this game, the strategy of [BKM19] selects an arbitrary  $h_k \in \mathcal{H}_{\varepsilon}(u_k)$  in round k. We will next show a strategy for selecting  $h_k$  that will allow for a faster convergence.

# 2.3 Warm-up: poly(log $n/\varepsilon^2$ ) Sample Complexity

So far, we reduced the hypothesis selection problem to solving the cutting-with-margin game. We next outline a solution for the cutting-with-margin game in  $\tilde{O}(\log n/\varepsilon^2)$  rounds. Since the implementation of each round requires  $O(\log n/\varepsilon^2)$  samples (see Section 2.1.1), this implies an algorithm for hypothesis selection with  $\tilde{O}(\log^2 n/\varepsilon^4)$  sample complexity.

First observe that an equivalent way of presenting the cutting-with-margin game lets the adversary pick in each round a halfspace  $H_k$  which is disjoint from the  $\ell_1$  ball of radius  $\varepsilon$  around  $h_k$ , and the game continues with  $\mathcal{H}_{k+1} = \mathcal{H}_k \cap H_k$ . This presentation is reminiscent of *Grunbaum's inequality* [Grü60], which guarantees that if the player picks the *centroid* (which is a standard way of defining the "center" of a body) of  $\mathcal{H}_k$  then  $vol(\mathcal{H}_{k+1}) \leq (1-e^{-1}) \cdot vol(\mathcal{H}_k)$ , where  $vol(\cdot)$  is the standard (Lebesgue) volume. While the centroid is an intuitive choice for our player, a counter strategy by the adversary will pick bodies that have small volumes but large diameters. Indeed, note that as long as the diameter of the body is greater than  $\varepsilon$ , the adversary can force at least one additional round. This shows that the volume is too crude of a measure for our game. Ideally, we would have wanted to use a different "centroid" that satisfies an analogous property with respect to the diameter (say,  $diameter(\mathcal{H}_{k+1}) \leq \frac{99}{100} \cdot diameter(\mathcal{H}_k)$ ). Unfortunately, no such object exists.

The approach we take for designing our player stems from the observation that if the player could always pick a point  $h_k \in \Delta_n$  that is close to the uniform distribution  $h^* = (\frac{1}{n}, \dots, \frac{1}{n})$ , then the game would have been solved in a few rounds. It is the easiest to see why when using the "primal" point of view from Section 2.1: indeed, assume  $u_k + \varepsilon \cdot 1_n$  is separated from  $\mathcal{P}$  by a hyperplane perpendicular to  $h^* = (\frac{1}{n}, \dots, \frac{1}{n})$ . Then, since  $u_{k+1} \geq u_k$  lies on the other side of that hyperplane, it follows that  $|u_{k+1} - u_k|_1 \geq \varepsilon n$ . So, when updating from  $u_k$  to  $u_{k+1}$ , the  $\ell_1$  norm increases by at least  $\varepsilon n$  (recall from Section 2.1 that in the [BKM19] strategy the  $\ell_1$  norm increases by only  $\Omega(\varepsilon)$  in each round). Thus, since in  $[0,1]^n$  the  $\ell_1$  norm is bounded by n, the total number of such steps is at most  $O(1/\varepsilon)$ . Of course, this strategy is impossible, as if  $h_1 = h^*$  then a ball of radius  $\varepsilon$  is disjoint from  $\mathcal{H}_k$ , for all k > 1.

Entropy as a progress measure. Inspired by the above intuition, our approach will be to set  $h_k \in \mathcal{H}_k$  to be as "close" to  $h^*$  as possible. Indeed, we select  $h_k \in \mathcal{H}_k$  that maximizes the entropy function (here we view the point  $h_k \in \Delta_n$  as a distribution). This corresponds to measuring the distance from the uniform distribution  $h^*$  using KL-divergence. The reason that the entropy function gives an efficient solution for our game boils down to that it is (i) strongly convex w.r.t  $\ell_1$  (as is evident by Pinsker's Inequality), (ii) bounded by  $\log(n)$  over the simplex. Roughly speaking, strong convexity means that in every step the entropy drops by  $\Omega(\varepsilon^2)$ . This, combined with the fact that the entropy is bounded by  $\log(n)$ , implies our  $\tilde{O}(\log(n)/\varepsilon^2)$  solution for the

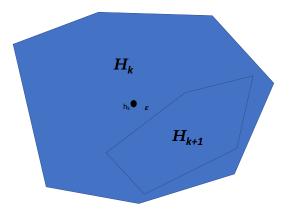


Figure 1: An illustration of the cutting-with-margin game: in each step k the player picks a point  $h \in \mathcal{H}_k$  and announces it to the adversary. The adversary then replies with  $\mathcal{H}_{k+1} \subseteq \mathcal{H}_k$  which is convex and disjoint from an  $\ell_1$  ball of radius  $\varepsilon$  around  $h_k$ . The players' goal is to empty the set as fast as possible (*i.e.*, to reach  $\mathcal{H}_k = \emptyset$ ), and the adversary's goal is to delay the player.

cutting-with-margin game<sup>14</sup>.

As discussed in the introduction, entropy and KL-divergence based strategies are often used in the context of optimization and regret minimization, basically for similar reasons (convexity and boundedness). However, our game is not defined by a cost function measuring the cost of each round separately, but rather, our "cost function" is the length of the game.

## 2.4 Near-Optimal Sample Complexity

In Section 2.3, we gave a hypothesis selection algorithm with  $\tilde{O}(\log^2 n/\varepsilon^4)$  samples, by solving the dual game. While this algorithm uses exponentially less samples than the one by [BKM19], it still sub-optimal. We next show how to obtain an algorithm with a near-optimal sample complexity of  $\tilde{O}(\log n/\varepsilon^2)$ , by first improving the dependence on n to  $\tilde{O}(\log n)$  (less involved), and then improving the dependence on  $\varepsilon$  to  $O(1/\varepsilon^2)$  (one of the main technical contributions of this paper). Since the sample complexity of our resulting algorithm (almost) matches Yatracos's, it can replace Yatracos's algorithm in density estimation algorithms to obtain a better approximation factor, while keeping the same low sample complexity.

<sup>&</sup>lt;sup>14</sup>Given that, it is natural to look for a strongly convex function over the simplex that is bounded by  $\ll \log(n)$ . However, no such function exists.

#### 2.4.1 Optimal Dependence on n

We revisit the basic observation from Section 2.1 that finding a distribution q satisfying  $(\forall i)$ :  $\mathsf{TV}(q,q_i) \leq \mathsf{TV}(p,q_i) + \varepsilon$  suffices in order to get a 2-approximation for hypothesis selection (see Equation (1)). We observe that it also suffices to find q that only satisfies  $\mathsf{TV}(q,q_{i^*}) \leq \mathsf{TV}(p,q_{i^*}) + \varepsilon$  (recall that  $i^*$  minimizes  $\mathsf{TV}(p,q_i)$ ) for exactly the same reason:  $\mathsf{TV}(q,p) \leq \mathsf{TV}(q,q_{i^*}) + \mathsf{TV}(q_{i^*},p) \leq 2\mathsf{TV}(q_{i^*},p) + \varepsilon$ . Thus, it suffices for our algorithm to maintain the invariant  $(u_k)_{i^*} \leq (v(p))_{i^*}$ , instead of  $u_k \leq v(p)$ . This suggests that we can relax Item 1 in Lemma 2 and only require  $(u_{k+1})_{i^*} \leq (v(p))_{i^*}$  (in addition to  $u_k \leq u_{k+1}$ ).

Due to the above, had we known  $i^*$ , we would only shoot for a good approximation (to within  $c \cdot \varepsilon$ ) of  $(u_{k+1})_{i^*}$ , which means that Lemma 2 can use only  $O(1/\varepsilon^2)$  samples (to get a good approximation of  $p(F_{i^*})$ ). But, we don't know the identity of  $i^*$ . The crucial observation here is that this does not matter. We can use the same  $O(1/\varepsilon^2)$  samples to evaluate each of the n statistical queries corresponding to each of the coordinates of  $u_{k+1}$ . Of course, since we are using too few samples, some of these coordinates will not be well approximated. However, it is likely that each one by itself will, and, in particular, this will be the case for  $(u_{k+1})_{i^*}$ . In other words, since we only care about  $(u_{k+1})_{i^*}$ , we no longer have to pay for a costly union bound over all n coordinates. (We also show that Item 2 in Lemma 2 still holds under this approximation using an averaging argument).

#### 2.4.2 Optimal Dependence on $\varepsilon$

Recall that in each step of the cutting-with-margin game, the player picks a point  $h_k \in \mathcal{H}_k$ , and the adversary sets  $\mathcal{H}_{k+1} \subseteq \mathcal{H}_k$  by cutting away an  $\ell_1$  ball of radius  $\varepsilon$  around  $h_k$ . The algorithm we have so far uses  $\Omega(\log n/\varepsilon^4)$  samples from p: every round uses  $\Theta(1/\varepsilon^2)$  samples and  $\max_{h \in \mathcal{H}_{\varepsilon}(u_k)} \{\mathbb{H}(h)\}$  drops by  $\Omega(\varepsilon^2)$  (recall that, to begin with, the entropy is at most  $\log n$  and we want it to drop to 0).

To reduce the sample complexity, we move away from this "static" type of algorithms and design a "dynamic" algorithm whose number of samples per round may vary (but, will never exceed  $\Omega(1/\varepsilon^2)$ ). The important property of the new algorithm is that if the algorithm samples more points from p, then the adversary cuts away a larger  $\ell_1$  ball around  $h_k$ . Specifically, if O(1) points are sampled then the radius of the removed ball is  $\varepsilon$ , and if  $O(1/\varepsilon^2)$  points are samples then the radius removed ball will be  $\Omega(1)$ . We will show that this coupling of the number of samples used in a step with the amount of progress made in that step (instead of using the maximum number of samples in every step and expecting the minimum progress) enables a win-win analysis which implies the desired saving in the sample complexity.

Bounding the radius of the removed ball. To explain how this idea is implemented, we need to dive into the details of the algorithm. Recall that the algorithm aims to find a point v such that  $v_{i^*} \leq \mathsf{TV}(p, q_{i^*}) + \varepsilon$ , and for which  $\mathcal{H}_{\varepsilon}(v) = \emptyset$ . Assume that the current point  $u_k$  satisfies  $d_{\infty}(u_k, \mathcal{P}) = d \gg \varepsilon$  (which means  $\mathcal{H}_d(u_k) = \emptyset$ ) and that we aim at reducing the distance to, say,  $\frac{3d}{4}$ . That is, we want to get to a point u such that  $d_{\infty}(u, \mathcal{P}) \leq \frac{3d}{4}$ , or, equivalently,  $\mathcal{H}_{\frac{3d}{4}}(u) = \emptyset$ . Recall from Section 2.1 that towards this, we pick a violated test  $h_k \in \mathcal{H}_{\frac{3d}{4}}(u_k)$  which, by applying Lemma 2, yields the new point  $u_{k+1} \in [0,1]^n$ . Of course, the lemma uses samples from p to compute this  $u_{k+1}$ . As we soon see, in some cases it will be worthwhile for our algorithm to only compute

a crude approximation of this  $u_{k+1}$  using fewer samples. Part of the difficulty is to decide on the quality of this approximation without knowing  $u_{k+1}$ .

Nevertheless, imagine for a moment that the algorithm does know this  $u_{k+1}$  and uses it as its next point. How much "progress" does this imply in the cutting-with-margin game? That is, how much smaller is  $\mathcal{H}_{\frac{3d}{4}}(u_{k+1})$  compared to  $\mathcal{H}_{\frac{3d}{4}}(u_k)$ ? Denote  $w_k = u_{k+1} - u_k$ . We next show that  $\mathcal{H}_{\frac{3d}{4}}(u_{k+1})$  is disjoint from an  $\ell_1$  ball of radius

$$r = \frac{d}{8\|w_k\|_{\infty}} \tag{3}$$

around  $h_k$  (we wish for r to be as large as possible). Intuitively, if  $||w_k||_{\infty}$  is small, it means that we have made progress in many coordinates (though the progress in each might be relatively small). Since we are getting close to  $\mathcal{P}$  in many directions, this should imply that  $u_{k+1}$  passes many of the tests  $h_k$  that were violated by  $u_k$ , and thus that  $\mathcal{H}_{\frac{3d}{2}}(u_{k+1})$  is much smaller.

More formally, let  $h \in \mathcal{H}_{\frac{3d}{4}}(u_{k+1})$ , Equation (3) follows from:

$$||h_k - h||_1 \cdot ||w_k||_{\infty} \ge (h_k - h) \cdot (u_{k+1} - u_k) \ge \frac{d}{8}.$$

Here, the first inequality is due Hölder's Inequality. The second inequality is because  $h_k \cdot (u_{k+1} - u_k) \ge \frac{3d}{8}$  (due to Lemma 2, Item 2) and because  $h \cdot (u_{k+1} - u_k) \le \frac{d}{4}$  (since  $h \in \mathcal{H}_{\frac{3d}{4}}(u_{k+1})$  it holds that  $h \cdot u_{k+1} + \frac{3d}{4} < \min_{v \in \mathcal{P}} h \cdot v$ , while since  $h \notin \mathcal{H}_d(u_k) = \emptyset$  it holds that  $h \cdot u_k + d \ge \min_{v \in \mathcal{P}} h \cdot v$ ).

Our "win-win" strategy. The take home message from the above discussion is that:

If 
$$||w_k||_{\infty}$$
 is small then  $\mathcal{H}_{\frac{3d}{4}}(u_{k+1})$  is small.

We next show that this relation leads us to a "win-win" situation: if  $||w_k||_{\infty}$  is large, it suffices to only crudely approximate  $w_k$ , and we save on samples. However, if  $||w_k||_{\infty}$  is small,  $\mathcal{H}_{\frac{3d}{4}}(u_{k+1})$  is small and we made a lot of progress towards ruling out all violated tests.

To see the relation between  $||w_k||_{\infty}$  and the number of samples required to approximate  $w_k$ , first assume that  $w_k$  is uniform over a set of coordinates of size m (i.e., for every  $i \in [n]$ , either  $(w_k)_i = 1/m$  or  $(w_k)_i = 0$ ). Now, if m is small than all non-zeros coordinates of  $w_k$  are large, and thus  $w_k$  can be reasonably approximated with few samples. (In fact, the number of samples scales with  $(1/||w_k||_{\infty})^2$ ).

**Slicing.** Of course,  $w_k$  may not be uniform on a set. To deal with such  $w_k$ 's, we partition  $w_k$  to  $\log(1/d)$  many "slices"  $w_k = w_k^1 + \ldots + w_k^{\log(1/d)}$  such that each  $w_k^\ell$  is almost uniform over a set (specifically, for  $\ell < \log(1/d)$ , each of the coordinates of  $w_k^\ell$  is either 0 or in  $(2^{-\ell}, 2^{-(\ell-1)}]$ ). We then try to identify a slice with a significant contribution to  $h_k \cdot w_k = \sum_{\ell \in [\log(1/d)]} h_k \cdot w_k^\ell$  (recall that  $h_k \cdot w_k \ge \frac{3d}{8}$  due to Lemma 2, Item 2). However, since  $w_k$  is not known to the algorithm, we use samples to learn it "slice-by-slice", starting by approximating  $w_k^1$ , the slice containing the largest values and requiring the least number of samples to estimate, and continuing to the slices that require more samples, until reaching a "good" slice. We mention that this slice-searching process is equivalent to playing the dual game with different  $\varepsilon$  values.

## 3 Preliminaries

#### 3.1 Notation

Let  $n \in \mathbb{N}$ . For  $u, v \in \mathbb{R}^n$ , we write  $u \geq v$  if  $(\forall i \in [n]) : u_i \geq v_i$ . We use  $u \cdot v := \sum_{i \in [n]} u_i v_i$  to denote the standard inner product of u and v.

For  $p \in [1, \infty]$ , we denote by  $\|\cdot\|_p$  the  $\ell_p$  norm. For  $u \in \mathbb{R}^n$  and  $r \geq 0$ , let  $B_p(u, r)$  denote a ball of radius r with respect to  $\ell_p$  that is centered at u,

$$B_p(u,\varepsilon) = \{ v \in \mathbb{R}^n : ||v - u||_p \le r \}.$$

Let  $\Delta_n$  denote the simplex of probability vectors in  $\mathbb{R}^n$ ,

$$\Delta_n := \{ h \in \mathbb{R}^n : \sum_{i \in [n]} h_i = 1, \ (\forall i) : h_i \ge 0 \}.$$

The entropy function is denoted by H and the Kullback-Leibler divergence by KL.

### 3.2 Definition of the Hypothesis Selection Problem

Let  $\mathcal{X}$  be a domain and let  $\Delta(\mathcal{X})$  denote the set of all probability distributions over  $\mathcal{X}$ . We assume that either (i)  $\mathcal{X}$  is finite in which case  $\Delta(\mathcal{X})$  is identified with the set of  $|\mathcal{X}|$ -dimensional probability vectors, or (ii)  $\mathcal{X} = \mathbb{R}^d$  in which case  $\Delta(\mathcal{X})$  is the set of Borel probability measures.

Let  $\mathcal{Q} \subseteq \Delta(\mathcal{X})$  be a set of distributions. We focus on the case where  $\mathcal{Q}$  is finite and denote its size by n. Let  $\alpha > 0$ , we say that  $\mathcal{Q}$  is  $\alpha$ -learnable with sample complexity  $m(n, \varepsilon, \delta)$  if there is a (possibly randomized) algorithm A such that for every  $\varepsilon, \delta > 0$  and every target distribution  $p \in \Delta(\mathcal{X})$ , if A receives as input at least  $m(n, \varepsilon, \delta)$  independent samples from p then it outputs a distribution q such that

$$\mathsf{TV}(p,q) \le \alpha \cdot \mathsf{opt} + \varepsilon,$$

with probability at least  $1 - \delta$ , where  $\mathsf{opt} = \min_{q \in \mathcal{Q}} \mathsf{TV}(p,q)$  and  $\mathsf{TV}(p,q) = \sup_{A \subseteq \mathcal{X}} \{p(A) - q(A)\}$  is the total variation distance. We say that  $\mathcal{Q}$  is properly  $\alpha$ -learnable if it is  $\alpha$ -learnable by a proper algorithm; namely an algorithm that always outputs  $q \in \mathcal{Q}$ .

**Distances vectors and sets.** Let  $Q = \{q_1, \ldots, q_n\} \subseteq \Delta(\mathcal{X})$ , and let p be a distribution. The TV-distance vector of p relative to the  $q_i$ 's is the vector  $v(p) = v_Q(p) = (\mathsf{TV}(p, q_i))_{i=1}^n$ .

Following [BKM19], our algorithm is based on the next claim which shows that in order to find q such that  $\mathsf{TV}(q,p) \leq 2 \min_i \mathsf{TV}(q_i,p) + \varepsilon$  it suffices to find q such that  $v(q) \leq v(p) + \varepsilon \cdot 1_n$ .

**Lemma 3.** Let q, p such that  $v(q) \le v(p) + \varepsilon \cdot 1_n$ . Then  $\mathsf{TV}(q, p) \le 2 \min_i \mathsf{TV}(q_i, p) + \varepsilon$ .

*Proof.* Follows directly by the triangle inequality; indeed, let  $q_i$  be a minimizer of  $\mathsf{TV}(\cdot, p)$  in  $\mathcal{Q}$ . Then,  $\mathsf{TV}(q, p) \leq \mathsf{TV}(q, q_i) + \mathsf{TV}(q_i, p) \leq (\mathsf{TV}(p, q_i) + \varepsilon) + \mathsf{TV}(q_i, p) = 2\mathsf{TV}(q_i, p) + \varepsilon$ .

Next, we explore which  $v \in \mathbb{R}^n$  are of the form v = v(p) for some  $p \in \Delta(\mathcal{X})$ . For this we make the following definition. A vector  $v \in \mathbb{R}^n$  is called a TV-distance dominating vector if  $v \geq v(p)$  for some distribution p. Define  $\mathcal{P}_{\mathcal{Q}}$  to be the set of all dominating distance vectors.

Claim 4.  $\mathcal{P}_{\mathcal{Q}}$  is convex and upward-closed<sup>15</sup>.

*Proof.* That  $\mathcal{P}_{\mathcal{Q}}$  is upward-closed is trivial. Convexity follows since  $\mathsf{TV}(\cdot, \cdot)$  is convex in both of its arguments.

### 3.3 Pythagorian Theorem for KL

We will use the following Pythagorian theorem for the KL divergence, the version here is taken from [PW15].

**Lemma 5.** Let  $\mathcal{X}$  be a set, let  $\mathcal{E} \subseteq \Delta(\mathcal{X})$  be a convex set of distributions, and let  $p \in \Delta(\mathcal{X})$  be a distribution. Let  $q^* = \arg\min_{q \in \mathcal{E}} \{ \mathsf{KL}(q, p) \}$ . Then, for all  $q \in \mathcal{E}$  it holds that

$$\mathsf{KL}(q, p) \ge \mathsf{KL}(q, q^*) + \mathsf{KL}(q^*, p).$$

*Proof.* If  $\mathsf{KL}(q,p) = \infty$ , then we are done. So, we can assume  $\mathsf{KL}(q,p) < \infty$ , which also implies that  $\mathsf{KL}(q^*,p) < \infty$ . For  $\theta \in [0,1]$ , form the convex combination  $q^{(\theta)} = (1-\theta)q^* + \theta q$ . Since  $q^*$  is the minimizer of  $\mathsf{KL}(q,p)$ , then

$$0 \le \left. \frac{\partial}{\partial \theta} \right|_{\theta=0} \mathsf{KL}(q^{(\theta)}, p) = \mathsf{KL}(q, p) - \mathsf{KL}(q, q^*) - \mathsf{KL}(q^*, p),$$

If we view the picture above in the Euclidean setting, the "triangle" formed by p,  $q^*$  and q (for  $q^*$ , q in a convex set, p is outside the set) is always obtuse, and is a right triangle only when the convex set has a "flat face". In this sense, the divergence is similar to the squared Euclidean distance, and the above theorem is sometimes called the Pythagorean theorem.

An assumption. Our analysis uses the Minimax Theorem for zero-sum games [vN28] for the same purpose that it was used in [BKM19]. Therefore, we will assume a setting (i.e., the domain  $\mathcal{X}$  and the class of distributions  $\mathcal{Q}$ ) in which this theorem is valid. Alternatively, one could state explicit assumptions such as finiteness of  $\mathcal{X}$  or forms of compactness under which it is known that the Minimax Theorem holds. However, we believe that the presentation benefits from avoiding such explicit technical assumptions and simply assuming the Minimax Theorem as an "axiom" in the discussed setting.

# 4 A Geometric Game from Hypothesis Selection

We next describe a geometric game, called the  $(\mathcal{P}, \varepsilon)$ -primal game. This game is between a player and an adversary, where  $\mathcal{P} \subseteq [0,1]^n$  is a given upwards-closed and nonempty convex body, and  $\varepsilon \geq 0$  is a margin parameter. Both  $\mathcal{P}$  and  $\varepsilon$  are known to both the player and the adversary. The game proceeds in rounds roughly as follows: the player starts at position  $u_0 = \vec{0} \in [0,1]^n$  and its goal is to get sufficiently close to  $\mathcal{P}$  as fast as possible. Let  $u_k$  denote the position of the player in

The Recall that upwards-closed means that whenever  $v \in \mathcal{Q}_{\mathcal{F}}$  and  $u \geq v$  then also  $u \in \mathcal{Q}_{\mathcal{F}}$ .

### The $(\mathcal{P}, \varepsilon)$ -Primal Game

Let  $\mathcal{P} \subseteq [0,1]^n$  be a nonempty convex set which is upward closed.

- 1. Set k = 0 and  $u_0 = \vec{0}$ .
- 2. While  $u_k + \varepsilon \cdot 1_n \notin \mathcal{P}$  (equivalently  $\mathcal{H}_{\mathcal{P},\varepsilon}(u_k) \neq \emptyset$ )
  - (a) The player picks a normal  $h_k \in \mathcal{H}_{\mathcal{P},\varepsilon}(u_k)$  to a hyperplane tangent to  $\mathcal{P}$  which separates  $u_k + \varepsilon \cdot 1_n$  from  $\mathcal{P}$ , and announces it to the adversary.
  - (b) The adversary replies with a point  $u_{k+1}$  whose every coordinate is at least as great as that of  $u_k$  and is  $\varepsilon/2$ -close to the hyperplane tangent to  $\mathcal{P}$  whose normal is  $h_k$ , *i.e.*,

$$u_{k+1} \ge u_k$$
 and  $h_k \cdot u_{k+1} \ge \min_{p \in \mathcal{P}} \{h_k \cdot p\} - \varepsilon/2.$  (4)

(c) Set k = k + 1.

Figure 2: The Primal Game.

round k; if  $u_k + \varepsilon \cdot 1_n \in \mathcal{P}$  then the player wins the game. Else, the player picks a tangent hyperplane to  $\mathcal{P}$  which separates  $u_k + \varepsilon \cdot 1_n$  from  $\mathcal{P}$  (such a hyperplane must exist since  $u_k + \varepsilon \cdot 1_n \notin \mathcal{P}$ ), announces it to the adversary, and the adversary picks the player's next position  $u_{k+1}$  to be any point such that  $u_{k+1} \geq u_k$  and  $u_{k+1}$  is  $\varepsilon/2$ -close to the tangent hyperplane chosen by the player. The  $(\mathcal{P}, \varepsilon)$ -primal game is formally described in Fig. 2. It uses the following notation:

$$\mathcal{H}_{\mathcal{P},\varepsilon}(u) = \left\{ h \in \Delta_n : h \cdot (u + \varepsilon \cdot 1_n) = h \cdot u + \varepsilon < \min_{p \in \mathcal{P}} \{h \cdot p\} \right\}.$$

In words,  $\mathcal{H}_{\mathcal{P},\varepsilon}(u)$  is the set of normals  $h \in \Delta_n$  to hyperplanes separating  $u + \varepsilon \cdot 1_n$  from  $\mathcal{P}$ . Note that the assumption  $h \in \Delta_n$  does not lose generality, because  $\mathcal{P}$  is upwards-closed and therefore for any  $u \in [0,1]^n$ ,  $u \notin \mathcal{P}$ , any hyperplane separating u and  $\mathcal{P}$  has a normal of this form. (See Claim 5 in [BKM19] for a proof of this fact.) Thus, by the hyperplane separation theorem,  $\mathcal{H}_{\mathcal{P},\varepsilon}(u_k) = \emptyset$  if and only if  $u_k + \varepsilon \cdot 1_n \in \mathcal{P}$ . Also observe that since  $\mathcal{P}$  is a convex, the set  $\mathcal{H}_{\mathcal{P},\varepsilon}(u)$  is convex for every  $u \in \mathbb{R}^n$ .

Winning Strategies. Let player be a strategy<sup>16</sup> for the player in the  $(\mathcal{P}, \varepsilon)$ -primal game. A sequence  $\vec{0} = u_0 \le u_1 \le ... \le u_t$  is a sequence of legal-adversary moves with respect to player if for every k < t,

- $\mathcal{H}_{\mathcal{P},\varepsilon}(u_k) \neq \emptyset$  and
- $h_k \cdot u_{k+1} \ge \min_{p \in \mathcal{P}} \{h_k \cdot p\} \varepsilon/2$ , where  $h_k = h_k(u_k; u_{< k}, h_{< k}) \in \mathcal{H}_{\mathcal{P}, \varepsilon}(u_k)$  is the normal picked by player in round k.

<sup>&</sup>lt;sup>16</sup>That is, in every round k, the strategy player provides a rule for picking  $h_k \in \mathcal{H}_{\mathcal{P},\varepsilon}(u_k)$ .

We say that the strategy player wins the  $(\mathcal{P}, \varepsilon)$ -primal game in at most r rounds if no adversary can force the game to last more than r rounds. That is, for every sequence  $\vec{0} = u_0 \le u_1 \le \ldots \le u_t$  of legal adversary-moves with respect to player,

$$u_t + \varepsilon \cdot 1_n \notin \mathcal{P} \implies t < r.$$

Similarly, let adv be a strategy<sup>17</sup> for the adversary in the  $(\mathcal{P}, \varepsilon)$ -primal game. A sequence  $h_0, \ldots, h_{t-1} \in \Delta_n$  is a sequence of legal-player moves with respect to adv if for every k < t,  $h_k \in \mathcal{H}_{\mathcal{P},\varepsilon}(u_k) \neq \emptyset$ , where  $u_k = u_k(h_{k-1}; u_{< k-1}, h_{< k-1}) \geq u_{k-1}$  is the point picked by adv in round k. We say that the strategy adv <u>forces the</u>  $(\mathcal{P}, \varepsilon)$ -primal game to last at least r rounds if for every sequence  $h_1, \ldots, h_{t-1} \in \Delta_n$  of legal-player moves with respect to adv,

$$u_t + 1_n \cdot \varepsilon \in \mathcal{P} \implies t \ge r.$$

## 4.1 Reducing Hypothesis Selection to the Primal Game

For all that follows, we fix a finite class of distributions  $Q = \{q_1, \ldots, q_n\}$  and  $\varepsilon > 0$ , and use the notation  $v(\cdot) = v_Q(\cdot)$ . We next show that if the  $(\mathcal{P}_Q, \varepsilon)$ -primal game is solvable in few rounds, then Q is 2-learnable with low sample complexity. The following lemma is implicitly proved in [BKM19]:

**Lemma 6.** If there exists a strategy player that wins the  $(\mathcal{P}_{\mathcal{Q}}, \varepsilon)$ -primal game in at most r rounds, then  $\mathcal{Q}$  is 2-learnable with sample complexity  $r'(\varepsilon, \delta) = O(r \cdot \frac{\log n + \log r + \log(1/\delta)}{\varepsilon^2})$ .

The reduction is described in Fig. 3. It is based on Lemma 3 and computes the output distribution q by finding  $v \in \mathcal{P}_{\mathcal{Q}}$  such that  $v \leq v(p) + \varepsilon \cdot n$ .

The following lemma is the crux of the reduction. It is used to show that the adversary induced by the algorithm is a valid adversary for the  $(\mathcal{P}_{\mathcal{Q}}, \varepsilon)$ -primal game, and provides a bound on the number of samples from p which are required to compute the adversary's move.

**Lemma 7.** Let  $p \in \Delta(\mathcal{X})$  and let  $\alpha, \beta > 0$ . Then, given  $m = O(\frac{\log n + \log(1/\beta)}{\alpha^2})$  independent samples from an unknown distribution p and  $h \in \Delta_n$  as an input, one can output a point  $z \in [0,1]^n$  that satisfies the following with probability  $\geq 1 - \beta$ :

1. 
$$h \cdot z \ge \min_{v \in \mathcal{P}_{\mathcal{Q}}} \{h \cdot v\} - \alpha$$
.

$$2. z \leq v(p).$$

In words, this lemma provides a procedure that, given a hyperplane tangent to  $\mathcal{P}_{\mathcal{Q}}$  and m samples from the target distribution p, outputs a point  $z \leq v(p)$  which is  $\alpha$ -close to the tangent.

<sup>&</sup>lt;sup>17</sup>That is, in every round k, the strategy adv provides a rule for picking  $u_{k+1}$  that satisfies Equation (4).

*Proof of Lemma 7.* By the Minmax Theorem [vN28]:

$$\min_{v \in \mathcal{P}_{\mathcal{Q}}} \{h \cdot v\} = \min_{p' \in \Delta(\mathcal{X})} \sum_{i \in [n]} h_i \cdot v(p')_i \qquad (By \text{ definition of } \mathcal{P}_{\mathcal{Q}}.)$$

$$= \min_{p' \in \Delta(\mathcal{X})} \sum_{i \in [n]} h_i \cdot \mathsf{TV}(p', q_i) \qquad (By \text{ definition of } v(\cdot))$$

$$= \min_{p' \in \Delta(\mathcal{X})} \sum_{i \in [n]} h_i \max_{f_i: \mathcal{X} \to [0, 1]} \{\mathbb{E}_{p'}[f_i] - \mathbb{E}_{q_i}[f_i]\} \qquad (By \text{ definition of } \mathsf{TV}(\cdot, \cdot).)$$

$$= \min_{p' \in \Delta(\mathcal{X})} \max_{f_i: \mathcal{X} \to [0, 1]} \sum_{i \in [n]} h_i(\mathbb{E}_{p'}[f_i] - \mathbb{E}_{q_i}[f_i])$$

$$= \max_{f_i: \mathcal{X} \to [0, 1]} \min_{p' \in \Delta(\mathcal{X})} \sum_{i \in [n]} h_i(\mathbb{E}_{p'}[f_i] - \mathbb{E}_{q_i}[f_i]). \qquad (By \text{ the Minmax Theorem.})$$

Let  $F_i$  for  $i \in [n]$  be maximizers of the last expression. That is,

$$(F_1,\ldots,F_n) = \operatorname{argmax}_{(f_1,\ldots,f_n)} \min_{p' \in \Delta(\mathcal{X})} \sum_{i \in [n]} h_i(\mathbb{E}_{p'}[f_i] - \mathbb{E}_{q_i}[f_i]).$$

By the above derivation:

$$\min_{p' \in \Delta(\mathcal{X})} \sum_{i \in [n]} h_i(\mathbb{E}_{p'}[F_i] - \mathbb{E}_{q_i}[F_i]) = \min_{v \in \mathcal{P}_{\mathcal{Q}}} \{h \cdot v\}.$$
 (5)

Note that the  $F_i$ 's depend only on the class  $\mathcal{Q}$  and the direction h; in particular they do not depend on p. Thus, the  $F_i$ 's can be computed by the algorithm. Define the point  $w \in [0,1]^n$  by

$$w_i = \mathbb{E}_p[F_i] - \mathbb{E}_{q_i}[F_i],$$

and observe that w can be approximated given samples from p. Note that w satisfies:

- 1.  $\sum_{i \in [n]} h_i w_i \ge \min_{v \in \mathcal{P}_{\mathcal{Q}}} \{h \cdot v\}$ . (By Equation (5).)
- 2.  $w_i = \mathbb{E}_p[F_i] \mathbb{E}_{q_i}[F_i] \le \max_{f_i: \mathcal{X} \to [0,1]} \{ \mathbb{E}_p[f_i] \mathbb{E}_{q_i}[f_i] \} = \mathsf{TV}(p, q_i) = v(p)_i.$

Thus, it suffices to output a point z such that  $w \ge z \ge w - \alpha \cdot 1_n$ . This can be done using the  $m = O(\frac{\log n + \log(1/\beta)}{\alpha^2})$  samples from p as follows: use the samples to approximate  $\mathbb{E}_p[F_i]$ . That is, let

$$\mathbb{E}_{\hat{p}}[F_i] = \frac{1}{m} \sum_{j=1}^m F_i(x_j),$$

where  $x_1, \ldots, x_m$  are the m independent samples drawn from p. By a Chernoff and union bounds, we have  $|\mathbb{E}_{\hat{p}}[F_i] - \mathbb{E}_p[F_i]| \leq \alpha/2$ , simultaneously for all  $i \leq n$ . Therefore, the estimates  $\hat{z}_i = \mathbb{E}_{\hat{p}}[F_i] - \mathbb{E}_{q_i}[F_i]$  satisfy  $\hat{z}_i \in (w_i - \frac{\alpha}{2}, w_i + \frac{\alpha}{2})$ . Then, the desired vector z can be taken to be  $z = \hat{w} - \frac{\alpha}{2} \cdot 1_n$ .

With Lemma 3, we are ready to prove Lemma 6 which shows how to use a black-box strategy for the player in the primal game to get a 2-approximation algorithm.

#### A Hypothesis Selection Algorithm from the Primal Game.

Define the set  $\mathcal{P}$  in the Primal Game to be  $\mathcal{P}_{\mathcal{Q}}$ . (The Primal Game is described in Fig. 2)

- 1. Set k = 0 and  $u_0 = \vec{0}$ .
- 2. While  $\mathcal{H}_{\mathcal{P}_{\mathcal{O}},\varepsilon}(u_k) \neq \emptyset$ 
  - (a) Run player to get  $h_k = h_k(u_k; u_{< k}, h_{< k}) \in \mathcal{H}_{\mathcal{P}_{\mathcal{O}}, \varepsilon}(u_k)$ .
  - (b) Let  $z_k$  be the point z promised by Lemma 7 applied with  $u = u_k$ ,  $h = h_k$ , p = p,  $\alpha = \varepsilon/2$  and  $\beta = \delta/r$ . Define  $u_{k+1}$  by setting  $(u_{k+1})_i = \max\{(z_k)_i, (u_k)_i\}$  for all  $i \in [n]$ .
  - (c) Set k = k + 1.
- 3. Output  $v = u_k + \varepsilon \cdot 1_n$ .

Figure 3: A Hypothesis Selection Algorithm from the Primal Game.

Proof of Lemma 6. Let player be a strategy for the player that wins the  $(\mathcal{P}_{\mathcal{Q}}, \varepsilon)$ -primal game in r rounds. We will show that  $\mathcal{Q}$  is 2-learnable with  $r' = O(r \cdot \frac{\log n + \log r + \log(1/\delta)}{\varepsilon^2})$  samples. Let  $p \in \Delta(\mathcal{X})$  be the target distribution. The approach we use for deriving the learning algorithm is based on Lemma 3 by which it suffices to find a distribution  $q \in \Delta(\mathcal{X})$  such that  $v(q) \leq v(p) + \varepsilon \cdot 1_n$ . Observe that if we find  $v \in \mathcal{P}_{\mathcal{Q}}$  such that  $v(q) \leq v(p) + \varepsilon \cdot 1_n$ , a distribution  $v(q) \leq v(p) + \varepsilon \cdot 1_n$  be found.

Consider the algorithm for computing such v which is depicted in Figure 3. The algorithm is based on an execution of the primal game, where the player runs the strategy player (see Item 2a) and the adversary moves are based on Lemma 7 (see Item 2b).

First note that Lemma 7 is applies with confidence parameter  $\beta = \delta/r$  and error parameter  $\varepsilon/2$ . This implies that: (i) the total number of samples used by the algorithm is r times the sample complexity bound stated in Lemma 7 with  $\alpha = \varepsilon/2, \beta = \delta/r$ , which yields the stated bound on r'. (ii) With probability at least  $1 - \delta$ , the points  $z_k$  satisfy the guarantee in Lemma 7 for all  $k \le r$ . In the remainder of the proof we condition on this event.

We next claim that the adversary strategy given in Item 2b provides a sequence of legal-adversary moves w.r.t player. To this end, we need to show that the point  $u_{k+1}$  satisfies  $u_{k+1} \ge u_k$  and  $h_k \cdot u_{k+1} \ge \min_{p \in \mathcal{P}} \{h_k \cdot p\} - \varepsilon/2$ . The former is obvious from the definition of  $u_{k+1}$  in Item 2b. The latter follows since

$$h_k \cdot u_{k+1} \ge h_k \cdot z_k \qquad (u_{k+1} \ge z_k, h_k \ge 0)$$
  
 
$$\ge \min_{p \in \mathcal{P}} \{ h_k \cdot p \} - \varepsilon/2 \qquad (\text{Lemma 7})$$

Thus, since player wins after at most r rounds, the while loop in Item 2 must terminate in round  $t \leq r$  and the output v satisfies  $v \in \mathcal{P}_{\mathcal{Q}}$ .

Finally, it remains to show that  $v \leq v(p) + \varepsilon \cdot 1_n$ . We show that  $u_k \leq v(p)$  for every k by induction on k (this implies  $v = u_t + \varepsilon \cdot 1_n \leq v(p) + \varepsilon \cdot 1_n$ ): For k = 0, it is clearly the case that

 $u_0 = \vec{0} \le v(p)$ . Assume that the claim holds for some k and prove it for k+1. By the second item in Lemma 7,  $z_k \le v(p)$ , and by the induction hypothesis,  $u_k \le v(p)$ . This implies that  $u_{k+1} \le v(p)$ .

## 5 A Dual Game: Cutting-With-Margin

One of the key steps in our solution is to adapt a dual point of view, where the separators/directions  $h_k$  are thought of as points in the dual space. We next describe a second geometric game, called the  $(\mathcal{H}, \varepsilon)$ -cutting-with-margin game (in short,  $(\mathcal{H}, \varepsilon)$ -cutting game), which can be seen as a manifestation of the primal game as seen in the dual space. This game too is between a player and an adversary, where  $\mathcal{H} \subseteq \Delta^n$  is a given convex body, and  $\varepsilon \geq 0$  is a margin parameter. Both  $\mathcal{H}$  and  $\varepsilon$  are known to both the player and the adversary.

The dual game proceeds in rounds roughly as follows: at the beginning, the universe is the set  $\mathcal{H}_0 = \mathcal{H}$ . In round k, the player chooses a point  $h_k \in \mathcal{H}_k$ . The adversary then restricts the universe to a set  $\mathcal{H}_{k+1}$ , which must to be a convex subset of  $\mathcal{H}_k$  that is disjoint from  $B_1(h_k, \varepsilon)$ , an  $\ell_1$  ball around  $h_k$  with radius  $\varepsilon$ . If the new universe  $\mathcal{H}_{k+1}$  is not empty, the game continues to the next round. Else, the game ends. A formal description of the dual game is given in Fig. 4.

Winning Strategies. Let player\* be a strategy<sup>18</sup> for the player in the  $(\mathcal{H}, \varepsilon)$ -cutting game. A sequence  $\mathcal{H} = \mathcal{H}_0 \supseteq \mathcal{H}_1 \supseteq \ldots \supseteq \mathcal{H}_t$  is a sequence of *legal-adversary moves with respect to* player\* if for every k < t,

- $\mathcal{H}_k \neq \emptyset$  and
- $\mathcal{H}_{k+1} \cap B_1(h_k, \varepsilon) = \emptyset$ , where  $h_k = h_k(\mathcal{H}_k; \mathcal{H}_{< k}, h_{< k}) \in \mathcal{H}_k$  is the point picked by player\* in round k.

We say that the strategy player\* wins the  $(\mathcal{P}, \varepsilon)$ -cutting game in at most r rounds if no adversary can force the game to last more than r rounds. That is, for every sequence  $\mathcal{H} = \mathcal{H}_0 \supseteq \mathcal{H}_1 \supseteq \ldots \supseteq \mathcal{H}_t$  of legal adversary-moves with respect to player\*,

$$\mathcal{H}_t \neq \emptyset \implies t < r.$$

Similarly, let  $\mathsf{adv}^*$  be a strategy<sup>19</sup> for the adversary in the  $(\mathcal{P}, \varepsilon)$ -primal game. A sequence  $h_0, \ldots, h_{t-1} \in \Delta_n$  is a sequence of legal-player moves with respect to  $\mathsf{adv}^*$  if  $h_k \in \mathcal{H}_k$  for every k < t, where  $\mathcal{H}_k = h_k(\mathcal{H}_{k-1}; \mathcal{H}_{< k-1}, h_{< k-1})$  is the set picked by  $\mathsf{adv}^*$  in round k-1. We say that the strategy  $\mathsf{adv}^*$  forces the  $(\mathcal{H}, \varepsilon)$ -cutting game to last at least r rounds if for every sequence  $h_1, \ldots, h_{t-1} \in \Delta_n$  of legal-player moves with respect to  $\mathsf{adv}^*$ ,

$$\mathcal{H}_t = \emptyset \implies t \ge r.$$

<sup>&</sup>lt;sup>18</sup>That is, in every round k, the strategy player provides a rule for picking  $h_k \in \mathcal{H}_k$ .

<sup>&</sup>lt;sup>19</sup>That is, in every round k, the strategy adv provides a rule for picking  $\mathcal{H}_{k+1}$  as in Item 2b.

### The $(\mathcal{H}, \varepsilon)$ -Cutting-With-Margin Game

Let  $\mathcal{H} \subseteq \Delta^n$  be convex.

- 1. Set k = 0 and  $\mathcal{H}_0 = \mathcal{H}$ .
- 2. While  $\mathcal{H}_k \neq \emptyset$ 
  - (a) The player picks a point  $h_k \in \mathcal{H}_k$  and announces it to the adversary.
  - (b) The adversary picks a convex set  $\mathcal{H}_{k+1} \subseteq \mathcal{H}_k$  such that  $\mathcal{H}_{k+1} \cap B_1(h_k, \varepsilon) = \emptyset$ . (Observe that  $\mathcal{H}_{k+1} = \emptyset$  always satisfies the above conditions.)
  - (c) Set k = k + 1.

Figure 4: The Dual Game.

### 5.1 Reduction from the Primal Game

For an upwards-closed convex set  $\mathcal{P} \subseteq [0,1]^n$  and  $\varepsilon > 0$ , let

$$\mathcal{H}_{\mathcal{P}} = \mathcal{H}_{\mathcal{P},\varepsilon}(\vec{0}) = \left\{ h \in \Delta_n : h \cdot \vec{0} + \varepsilon = \varepsilon < \min_{u \in \mathcal{P}} \{ h \cdot u \} \right\}.$$

We next show that the round complexity of the  $(\mathcal{P}, \varepsilon)$ -primal game is at most the round complexity of the  $(\mathcal{H}_{\mathcal{P}}, \varepsilon)$ -cutting game.

**Lemma 8.** Let  $\mathcal{P} \subseteq [0,1]^n$  be an upwards-closed convex set and let  $\varepsilon \geq 0$ . If there exists a strategy player\* that wins the  $(\mathcal{H}_{\mathcal{P}}, \frac{\varepsilon}{4})$ -cutting game in at most r rounds, then there is a strategy player that wins the  $(\mathcal{P}, \varepsilon)$ -primal game in at most r rounds.

The proof of Lemma 8 uses the following lemma. Recall that  $\mathcal{H}_{\mathcal{P},\varepsilon}(u') = \{h \in \Delta_n : h \cdot u' + \varepsilon < \min_{u \in \mathcal{P}} h \cdot u\}.$ 

**Lemma 9.** Let  $\mathcal{P} \subseteq [0,1]^n$  be an upwards-closed convex set and let  $h \in \Delta_n$ ,  $u \in [0,1]^n$ . Then,

$$h \notin \mathcal{H}_{\mathcal{P},\frac{\varepsilon}{2}}(u) \implies \mathcal{H}_{\mathcal{P},\varepsilon}(u) \cap B_1(h,\frac{\varepsilon}{4}) = \emptyset.$$

In other words if  $h \cdot u \ge \min_{p \in \mathcal{P}} \{h \cdot p\} - \frac{\varepsilon}{2}$  then,  $\mathcal{H}_{\mathcal{P},\varepsilon}(u) \cap B_1(h, \frac{\varepsilon}{4}) = \emptyset$ .

*Proof.* Let  $h \notin \mathcal{H}_{\mathcal{P},\frac{\varepsilon}{2}}(u)$ , we need to show that  $\mathcal{H}_{\mathcal{P},\varepsilon}(u) \cap B_1(h,\frac{\varepsilon}{4}) = \emptyset$ . Define  $G: \Delta_n \to \mathbb{R}$  by

$$G(h) = \min_{p \in \mathcal{P}} \{h \cdot p\} - h \cdot u.$$

Thus, G(h) measures the distance between u and the hyperplane tangent to  $\mathcal{P}$  with normal h. Observe that for every  $\varepsilon' \geq 0$ :

$$\mathcal{H}_{\mathcal{P},\varepsilon'}(u) = \{ h \in \Delta_n : G(h) > \varepsilon' \}. \tag{6}$$

Note that

- (i)  $G(h) \leq \frac{\varepsilon}{2}$  (by Equation (6), because  $h \notin \mathcal{H}_{\mathcal{P},\frac{\varepsilon}{2}}(u)$ ).
- (ii) G is 1-Lipschitz with respect to  $\ell_1$ : Let  $h', h'' \in \Delta_n$  and let  $p^* := \arg\min_{p \in \mathcal{P}} \{h'' \cdot p\}$ . It holds that

$$\begin{split} G(h') - G(h'') &= \min_{p \in \mathcal{P}} \{h' \cdot p\} - \min_{p \in \mathcal{P}} \{h'' \cdot p\} - (h' - h'') \cdot u \\ &\leq h' \cdot p^* - h'' \cdot p^* - (h' - h'') \cdot u \\ &= (h' - h'') \cdot (p^* - u) \\ &\leq \|h' - h''\|_1 \cdot \|p^* - u\|_{\infty} \quad \text{(H\"older's inequality, } (\forall v, v') : v \cdot v' \leq \|v\|_1 \|v'\|_{\infty}) \\ &\leq \|h' - h''\|_1. \qquad \qquad (\|p^*\|_{\infty} \leq 1, \ \|u\|_{\infty} \leq 1, \ p^*, u \geq 0) \end{split}$$

Let  $h' \in B_1(h, \frac{\varepsilon}{4})$ , so  $||h' - h||_1 \leq \frac{\varepsilon}{4}$ . Thus,

$$G(h') \le G(h) + \frac{\varepsilon}{4}$$
 (By Item (ii).)  
  $\le \varepsilon$  (By Item (i).)

Thus, by Equation (6)  $h' \notin \mathcal{H}_{\mathcal{P},\varepsilon}(u)$ , and  $\mathcal{H}_{\mathcal{P},\varepsilon}(u) \cap B_1(h,\frac{\varepsilon}{2}) = \emptyset$ , as required.

Proof of Lemma 8. Let player\*( $\mathcal{H}$ ) be a strategy for the player that solves the  $(\mathcal{H}_{\mathcal{P}}, \frac{\varepsilon}{4})$ -cutting game in r rounds. Consider the reduction described in Figure 5 and the strategy player for the  $(\mathcal{P}, \varepsilon)$ -primal game which is described in Item 2b. Our goal is to show that player wins the game in at most r rounds. That is, let  $\vec{0} = u_0, \ldots, u_t$  be a sequence of legal-adversary moves w.r.t player such that  $u_t + \varepsilon \cdot 1_n \notin \mathcal{P}$ . We need to show that t < r. To this end it suffices to show that the sequence  $\{\mathcal{H}_k\}_{k=0}^t$ , defined in Item 2a is a sequence of legal-adversary moves w.r.t player\* and that  $\mathcal{H}_t \neq \emptyset$ . Indeed, by definition  $\mathcal{H}_0 = \mathcal{H}_{\mathcal{P},\varepsilon}(\vec{0}) = \mathcal{H}_{\mathcal{P},\varepsilon}$ . Next,

$$\mathcal{H}_{k+1} = \mathcal{H}_{\mathcal{P},\varepsilon}(u_{k+1}) \subseteq \mathcal{H}_{\mathcal{P},\varepsilon}(u_k) = \mathcal{H}_k,$$

because  $u_{k+1} \geq u_k$  and  $\mathcal{H}_{k+1}, \mathcal{H}_k \subseteq \Delta_n$  contains only nonnegative vectors. The last property we need to show in order to establish that the  $\mathcal{H}_k$ 's form a sequence of legal-adversary moves is that  $\mathcal{H}_{k+1} \cap B_1(h_k, \varepsilon) = \emptyset$  for k < t, which follows from Lemma 9 because  $h_k \cdot u_{k+1} \geq \min_{p \in \mathcal{P}} \{h_k \cdot p\} - \varepsilon/2$  (i.e.,  $h_k \notin \mathcal{H}_{\mathcal{P}, \frac{\varepsilon}{2}}$ ). Finally, it remains to show that  $\mathcal{H}_t \neq \emptyset$ . Indeed, by assumption,  $u_t + \varepsilon \cdot 1_n \notin \mathcal{P}$  and therefore there must be a hyperplane separating  $u_t + \varepsilon \cdot 1_n$  from  $\mathcal{P}$ . Hence, the normal to this hyperplane (normalized so that it is in  $\Delta_n$ ) belongs to  $\mathcal{H}_{\mathcal{P},\varepsilon}(u_t) = \mathcal{H}_t$  and witnesses  $\mathcal{H}_t \neq \emptyset$ .

5.2 Solution for the Cutting-With-Margin Game

**Theorem 10.** For every convex set  $\mathcal{H}$  and  $\varepsilon > 0$ , the  $(\mathcal{H}, \varepsilon)$ -cutting game is solvable in  $O\left(\frac{\log(n)}{\varepsilon^2}\right)$  rounds.

### Dual Player Strategy $\implies$ Primal Player Strategy.

Let  $\mathcal{P} \subseteq [0,1]^n$  be an upward-closed convex set, let  $\varepsilon > 0$ , and consider the  $(\mathcal{P}, \varepsilon)$ -primal Game. Let  $\mathsf{player}^*(\cdot)$  be a strategy for the player in the  $(\mathcal{H}_{\mathcal{P},\varepsilon}, \varepsilon/4)$ -cutting game.

- 1. Set k = 0 and  $u_0 = \vec{0}$ .
- 2. While  $u_k + \varepsilon \cdot 1_n \notin \mathcal{P}$  (equivalently  $\mathcal{H}_{\mathcal{P},\varepsilon}(u_k) \neq \emptyset$ )
  - (a) Let  $\mathcal{H}_k := \mathcal{H}_{\mathcal{P},\varepsilon}(u_k)$ .
  - (b) Run player\* to get  $h_k = h_k(\mathcal{H}_k; \mathcal{H}_{\leq k}, h_{\leq k}) \in \mathcal{H}_{\mathcal{P}, \varepsilon}(u_k)$  and announce  $h_k$  to the adversary.
  - (c) Let  $u_{k+1}$  denote the next point which is picked by the adversary. I.e.,  $u_{k+1} \ge u_k$  and  $h_k \cdot u_{k+1} \ge \min_{p \in \mathcal{P}} \{h_k \cdot p\} \varepsilon/2$ .
  - (d) Set k = k + 1.

Figure 5: A reduction which uses a black box access to a strategy player\* in the  $(\mathcal{H}_{\mathcal{P},\varepsilon},\varepsilon/4)$ -cutting game and produces a strategy player for the  $(\mathcal{P},\varepsilon)$ -primal game. This reduction is used in the proof of Lemma 8.

#### A Strategy for the Player in the Cutting-With-Margin Game.

In each round k, given a universe  $\mathcal{H}_k$  from the adversary, the player outputs  $h_k = \arg\max_{h' \in \mathcal{H}_k} \{\mathbb{H}(h')\}.$ 

Figure 6: A strategy for the player in the  $(\mathcal{H}, \varepsilon)$ -cutting game. Recall that  $\mathbb{H}(\cdot)$  denotes the entropy function, and that  $\mathcal{H} \subseteq \Delta_n$  and therefore  $\mathbb{H}(\cdot)$  is defined on every  $h \in \mathcal{H}$ .

*Proof.* Consider the strategy player\* for the player in the  $(\mathcal{H}, \varepsilon)$ -cutting game that is depicted in Figure 6.

Fix a sequence  $\mathcal{H} = \mathcal{H}_0 \supseteq \ldots \supseteq \mathcal{H}_t$  of legal-adversary moves w.r.t player\* such that  $\mathcal{H}_t \neq \emptyset$ . Our goal is to prove that  $t \leq O(\log n/\varepsilon^2)$ .

For  $k \leq t$ , let  $h_k = \arg \max_{h \in \mathcal{H}_k} \mathbb{H}(h)$  denote the point chosen by player\*. Note that  $h_k$  is well defined since  $\mathcal{H}_k \neq \emptyset$  for  $k \leq t$ . We next prove that for every k < r,

$$\mathbb{H}(h_k) - \mathbb{H}(h_{k+1}) \ge \frac{\varepsilon^2}{8}.\tag{7}$$

This implies the desired bound  $t \leq O\left(\frac{\log(n)}{\varepsilon^2}\right)$  as follows: since the entropy function satisfies  $0 \leq \mathbb{H}(h) \leq \log(n)$  for every  $h \in \Delta_n$ . In particular, by Eq. (7),

$$0 \le \mathbb{H}(h_t) \le \mathbb{H}(h_0) - \frac{t \cdot \varepsilon^2}{8} \le \log(n) - \frac{t \cdot \varepsilon^2}{8},$$

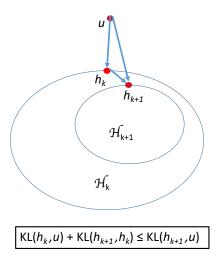


Figure 7: An illustration of the Pythagorean Theorem for Kullback-Leibler divergence as it is used in the proof of Theorem 10.

and therefore  $t \leq \frac{8\log(n)}{\varepsilon^2}$ , as required.

It remains to prove Equation (7). Since  $\mathcal{H} = \mathcal{H}_0 \supseteq \ldots \supseteq \mathcal{H}_t$  is a sequence of legal-adversary moves, it holds that  $\mathcal{H}_{k+1} \cap B_1(h_k, \frac{\varepsilon}{2}) = \emptyset$ . Since  $h_{k+1} \in \mathcal{H}_{k+1}$ , it holds that  $h_{k+1} \notin B_1(h_k, \frac{\varepsilon}{2})$ , and therefore,  $||h_{k+1} - h_k||_1 \ge \frac{\varepsilon}{2}$ . Let  $u \in \Delta_n$  denote the uniform distribution  $u = (\frac{1}{n}, \ldots, \frac{1}{n})$ .

$$\begin{split} \mathbb{H}(h_k) - \mathbb{H}(h_{k+1}) &= \mathsf{KL}(h_{k+1}, u) - \mathsf{KL}(h_k, u) \\ &\geq \mathsf{KL}(h_{k+1}, h_k) & \text{(Lemma 5, see reasoning below)} \\ &\geq \frac{1}{2} \|h_{k+1} - h_k\|_1^2 & \text{(Pinsker's Inequality)} \\ &\geq \frac{\varepsilon^2}{8}. & \text{(} \|h_{k+1} - h_k\|_1 \geq \frac{\varepsilon}{2}) \end{split}$$

For the second transition, the inequality  $\mathsf{KL}(h_{k+1}, u) - \mathsf{KL}(h_k, u) \geq \mathsf{KL}(h_{k+1}, h_k)$  follows from the Pythagorean Theorem for the Kullback-Leibler divergence, Lemma 5, by taking  $\mathcal{E} = \mathcal{H}_k$ , p = u,  $q^* = h_k$ ,  $q = h_{k+1}$ . For the theorem to apply, we need to use the facts that  $\mathcal{H}_k$  is convex, that  $h_k \in \mathcal{H}_k$ , and that  $h_{k+1} \in \mathcal{H}_{k+1} \subseteq \mathcal{H}_k$ . We also need

$$h_k = \arg \max_{h \in \mathcal{H}_k} \{ \mathbb{H}(h) \} = \arg \min_{h \in \mathcal{H}_k} \{ \mathsf{KL}(h, u) \}.$$

See Figure 7 for an illustration.

## 5.3 Hypothesis Selection with $poly(log n/\epsilon)$ Samples

The results obtained so far already suffice for constructing a hypothesis selection algorithm with sample complexity of  $poly(\log n/\varepsilon)$ , as suggested by the proposition below. The algorithm we obtain in this subsection will be refined in Section 6 to obtain an algorithm with near optimal sample complexity.

**Proposition 11.** Let Q be a finite class of distributions and let n = |Q|. Then, Q is 2-learnable with sample complexity

$$m(n, \varepsilon, \delta) = O\left(\frac{\log^2 n + \log n \log(1/\varepsilon) + \log n \log(1/\delta)}{\varepsilon^4}\right).$$

*Proof.* Let  $\mathcal{P} = \mathcal{P}_{\mathcal{Q}}$  denote the set of all dominating-distance vectors w.r.t  $\mathcal{Q}$ . By Lemma 6, it suffices to give a strategy for the player that wins the  $(\mathcal{P}, \varepsilon)$ -primal game in at most  $r = O(\log n/\varepsilon^2)$  rounds. The existence of such a strategy follows from Theorem 10, which yields a strategy for the player that wins the  $(\mathcal{H}_{\mathcal{P},\varepsilon}, \varepsilon/4)$  in  $O(\log n/\varepsilon^2)$  rounds, and by Lemma 8 which transforms this strategy to a strategy that wins the  $(\mathcal{P}, \varepsilon)$  game in the same number of rounds.

See Fig. 8 for a pseudo-code of the algorithm obtained by this series of reductions.  $\Box$ 

#### A 2-Approximation Algorithm for Hypothesis Selection with poly $(\log n/\varepsilon)$ Samples

Given: A class  $Q = \{q_1, \ldots, q_n\}$ , and a sampling access to a target distribution p and  $\varepsilon, \delta > 0$ . Output: A distribution  $p_0$  such that  $\mathsf{TV}(p_0, p) \leq 2 \min_i \mathsf{TV}(q_i, p) + \varepsilon$  with probability at least  $1 - \delta$ .

- 1. Let  $v^* = v(p) = (\mathsf{TV}(p, q_i))_i \in \mathbb{R}^n$ , and set  $u_0 = (0, \dots, 0) \in \mathbb{R}^n$ . (Note that  $v^*$  is not known to the algorithm)
- 2. For k = 1, ...
  - (a) If  $u_k + \varepsilon \cdot 1_n \in \mathcal{P} = \mathcal{P}_{\mathcal{Q}}$  then find q such that  $v(q) \leq u_k + \varepsilon \cdot 1_n$  and output it.
  - (b) Else, pick a separator  $h_k = \arg \max_{h \in \mathcal{H}_{\mathcal{P},\varepsilon}(x_k)} \{ \mathbb{H}(h) \}$  with maximum entropy.
  - (c) Draw  $m = O(\frac{\log n + \log \log n + \log(1/\varepsilon) + \log(1/\delta)}{\varepsilon^2})$  samples from p to compute  $u_{k+1}$  such that  $u_k \le u_{k+1} \le v^*$ , and

$$h_k \cdot u_{k+1} \ge \min_{u \in \mathcal{D}} \{h_k \cdot u\} - \frac{\varepsilon}{2}.$$

(See Lemma 7 for the computation of  $u_{k+1}$ .)

(d) Continue to the next iteration.

Figure 8: The hypothesis selection algorithm obtained by the reductions to the two games.

# 6 Obtaining Near Optimal Sample Complexity

In this section we prove Theorem 1, giving a 2-approximation algorithm for hypothesis selection with sample complexity that is tight up to lower-order terms. To this end, we first study a refined version of the primal game from Section 4, and then study a refined version of the hypothesis selection algorithm given in Figure 8 (that was based on our solution to the cutting-with-margin game).

#### The Refined $(\mathcal{P}, \varepsilon)$ -Primal Game

Let  $\mathcal{P} \subseteq [0,1]^n$  be a nonempty convex set which is upward closed,  $C_0$  here is a large constant.

- 1. Set k = 0,  $u_0 = \vec{0}$ , d = 1,  $d' = d/(C_0 \log(1 + 1/d))$ .
- 2. While  $d > \varepsilon/2$ 
  - (a) While  $\mathcal{H}_{\mathcal{P},d-d'}(u_k) \neq \emptyset$ 
    - i. The player picks a normal

$$h_k = \operatorname{argmax}_{h' \in \mathcal{H}_{\mathcal{P}, d-d'}(u_k)} \{ \mathbb{H}(h') \} \in \mathcal{H}_{\mathcal{P}, d-d'}(u_k)$$

to a hyperplane tangent to  $\mathcal{P}$  which separates  $u_k + (d - d') \cdot 1_n$  from  $\mathcal{P}$ .

ii. Run Refined Hypothesis Select algorithm<sup>a</sup> to reply with a point  $u_{k+1}$ , and an integer  $j \in \{0, \ldots, 2 + \lceil \log(1 + 1/d) \rceil \}$  which satisfy

$$u_{k+1} \ge u_k$$
 and  $\sum_{i} \min(2^{-j}, u_{k+1,i} - u_{k,i}) \cdot h_{k,i} > 2d'$  (8)

iii. Set k = k + 1.

- (b) Set d = d d',  $d' = d/(C_0 \log(1 + 1/d))$ .
- 3. Output a distribution r satisfying  $\mathsf{TV}(r, q_i) \leq u_{k,i} + d$  for all i.

Figure 9: The Refined Primal Game.

## 6.1 Refining the Primal Game

Consider the Refined Primal Game algorithm given in Figure 9. This algorithm uses the Refined Hypothesis Select algorithm that can be found in Figure 10. Let us analyze the Refined Primal Algorithm before presenting the Refined Hypothesis Select component. The critical property of the refined hypothesis select part will be that the number of samples needed to get (8) to hold with a given j is  $\tilde{O}(2^{2j})$ , which can be substantially smaller than  $1/\varepsilon^2$  required by the analogous step

<sup>&</sup>lt;sup>a</sup>Note that stage requires  $\tilde{O}(2^{2j})$  qureies

of the original algorithm. At the same time, when j is large (and thus the sample complexity cost is large), we get a stronger estimate of  $\Omega(2^j \cdot d')$  on the distance  $||h_k - h_{k+1}||_1$ , which translates into more progress towards reducing the value of  $\mathbb{H}(h_{k+1})$ , helping the algorithm terminate faster. Thus we get a win-win situation, where small j means fewer samples needed, and a large j means a lot of progress towards completion.

Claim 12. (Refined  $(\mathcal{P}, \varepsilon)$ -Primal Game – running time) Suppose that conditioned on outputting j, Refined Hypothesis Select terminates after  $\leq A \cdot 2^{2j}$  samples in expectation. Then the expected number of samples needed by the Refined  $(\mathcal{P}, \varepsilon)$ -Primal Game is  $\tilde{O}((A \cdot \log n)/\varepsilon^2)$ .

*Proof.* Note that the outer loop (Item 2), where d gets reduced runs a total of  $O(\log^2(1/\varepsilon))$  times, and therefore it suffices to analyze one execution of the loop to show that as long as  $d = \Omega(\varepsilon)$ , the number of samples used in reducing d to d - d' is bounded by  $\tilde{O}((A \cdot \log n)/\varepsilon^2)$ .

Consider a single iteration of the inner loop, we would like to lower bound the difference  $\mathbb{H}(h_k) - \mathbb{H}(h_{k+1})$ . By the exection of the algorithm  $\mathcal{H}_{\mathcal{P},d}(u_k) = \emptyset$ , and thus  $h_{k+1} \notin \mathcal{H}_{\mathcal{P},d}(u_k)$ . Therefore, by definition of  $\mathcal{H}_{\mathcal{P},d}(u_k)$  it holds:

$$\min_{v \in \mathcal{P}} \{ h_{k+1} \cdot v \} \le h_{k+1} \cdot u_k + d. \tag{9}$$

On the other hand,  $h_{k+1} \in \mathcal{H}_{\mathcal{P},d-d'}(u_{k+1})$ , and thus

$$\min_{v \in \mathcal{P}} \{ h_{k+1} \cdot v \} > h_{k+1} \cdot u_{k+1} + d - d'. \tag{10}$$

Putting equations (9) and (10) together, we get:

$$h_{k+1} \cdot u_{k+1} + d - d' < \min_{v \in \mathcal{P}} \{ h_{k+1} \cdot v \} \le h_{k+1} \cdot u_k + d.$$

Thus from above,

$$d' > h_{k+1} \cdot (u_{k+1} - u_k) \ge \sum_{i} \min(2^{-j}, u_{k+1,i} - u_{k,i}) \cdot h_{k+1,i} =$$

$$\sum_{i} \min(2^{-j}, u_{k+1,i} - u_{k,i}) \cdot h_{k,i} + \sum_{i} \min(2^{-j}, u_{k+1,i} - u_{k,i}) \cdot (h_{k+1,i} - h_{k,i}).$$

Applying (8) on the RHS we get:

$$d' > 2d' + \sum_{i} \min(2^{-j}, u_{k+1,i} - u_{k,i}) \cdot (h_{k+1,i} - h_{k,i}) > 2d' - 2^{-j} \cdot ||h_{k+1} - h_{k}||_{1}.$$

Therefore,

$$||h_{k+1} - h_k||_1 > d' \cdot 2^j. \tag{11}$$

By the same derivation as in the proof of Theorem 10, Equation (11) implies

$$\mathbb{H}(h_k) - \mathbb{H}(h_{k+1}) = \mathsf{KL}(h_{k+1}, u) - \mathsf{KL}(h_k, u)$$

$$\geq \mathsf{KL}(h_{k+1}, h_k) \qquad \qquad \text{(Lemma 5)}$$

$$\geq \frac{1}{2} \|h_{k+1} - h_k\|_1^2 \qquad \qquad \text{(Pinsker's Inequality)}$$

$$> \frac{d'^2}{2} \cdot 2^{2j} \qquad \qquad \text{(Equation (11))}$$

$$> \frac{\varepsilon^2 / \log^2(1/\varepsilon)}{4} \cdot 2^{2j} = \tilde{\Omega}(\varepsilon^2 \cdot 2^{2j}).$$

At the beginning of execution with a given d,  $\mathbb{H}(h_k) \leq \log n$ , and at the end it is at least 0. Each step causing a reduction by  $\tilde{\Omega}(\varepsilon^2 \cdot 2^{2j})$  takes  $\leq A \cdot 2^{2j}$  queries. Thus the total number of queries for a given value of d is bounded by  $\tilde{O}((A \log n)/\varepsilon^2)$ .

Claim 13. (Refined  $(\mathcal{P}, \varepsilon)$ -Primal Game – correctness) Let  $i \in [n]$  be any fixed index, which may depend on p and the q's but not on the execution of the algorithm. Suppose that at every step k of Refined Hypothesis Select, the probability

$$\Pr[u_{k+1,i} > \mathsf{TV}(p,q_i)] < \tilde{o}(\delta \varepsilon^2 / \log n).$$

Then the probability that  $\mathsf{TV}(r,q_i) > \mathsf{TV}(p,q_i) + \varepsilon/2$  is at most  $\delta$ .

*Proof.* At each step,  $\mathbb{H}(h_k)$  decreases by at least  $\tilde{\Omega}(\varepsilon^2)$ , and thus the total number of calls to the Refined Hypothesis Select algorithm is  $\tilde{O}((\log n)/\varepsilon^2)$ . Therefore, by union bound, except with probability  $<\delta$ , at each step k,  $u_{k,i} \leq \mathsf{TV}(p,q_i)$ . Therefore, the distribution r the algorithm outputs satisfies

$$\mathsf{TV}(r,q_i) \leq u_{k_{end},i} + d \leq \mathsf{TV}(p,q_i) + d < \mathsf{TV}(p,q_i) + \varepsilon/2.$$

## 6.2 Refining the Hypothesis Selection Algorithm

We next turn our attention to the Refined Hypothesis Select algorithm in Figure 10.

Note that the number of samples used by the Refined Hypothesis Select algorithm is spelled out explicitly. Therefore, our only task is to show that its success guarantees hold. Properties (ii) and (iii) holds due to stopping conditions of the algorithm. Next claim proves that Property (iv) holds.

Claim 14. Fix an index i. Assuming Refined Hypothesis Select algorithm does not output 'Fail', the probability of the event

$$\Pr[(v_i > u_i) \land (\mathbb{E}_p[F_i] - \mathbb{E}_{q_i}[F_i] \le v_i + 2^{-j-2})] < \gamma.$$
(15)

*Proof.* Note that  $v_i = \max(u_i, w_{ji} - 2^{-j-1})$ . Therefore,  $v_i > u_i$  iff  $w_{ji} > u_i + 2^{-j-1}$ . Recall that  $w_{ji} = \sum_{k=1}^{m_j} F_i(x_k) - \mathbb{E}_{q_i}[F_i]$ . Therefore event  $\{v_i > u_i\}$  dominated by the event

### The Refined Hypothesis Select Algorithm

Given  $d, d' = d/(C_0 \log(1+1/d))$ , error parameter  $\gamma > 0$ , a point u such that  $\mathcal{H}_{\mathcal{P}_{\mathcal{Q}},d}(u) = \emptyset$ , and a distribution  $h \in \mathcal{H}_{\mathcal{P}_{\mathcal{Q}},d-d'}(u)$  the algorithm will output  $j \in \{0,\ldots,2+\lceil \log(1+1/d)\rceil\}$ , a point v and n functions  $F_i: \mathcal{X} \to [0,1], i=1,\ldots,n$  such that the following properties hold:

- (i) The algorithm outputs 'success' with probability  $> 1 \gamma$ , where the failure event only depends on the randomness of the samples the algorithm receives;
- (ii)  $v \geq u$ ;
- (iii)  $\sum_{i} \min(2^{-j}, v_i u_i) \cdot h_i > 2d'$
- (iv) For any  $i \in [n]$  which is fixed in advance (unknown to the algorithm) if  $v_i > u_i$ , then except with probability  $\gamma$ ,  $\mathbb{E}_p[F_i] \mathbb{E}_{q_i}[F_i] > v_i + 2^{-j-2}$ .

Algorithm:

1. Let  $\{F_i\}_{i=1}^n$  be as in the proof of Lemma 7:  $F_i: \mathcal{X} \to [0,1]$  such that

$$\sum_{i \in [n]} h_i \cdot (\mathbb{E}_p[F_i] - \mathbb{E}_{q_i}[F_i]) \ge \min_{v \in \mathcal{P}_{\mathcal{Q}}} \{h \cdot v\}.$$
 (12)

- 2. For  $j \in \{0, \dots, 2 + \lceil \log(1 + 1/d) \rceil \}$ :
  - (a) Use  $m_j := C_1 \log(\log(1/d)/\gamma) \cdot 2^{2j}$  samples  $\{x_k\}_{k=1}^{m_j}$  from P to generate empirical estimates

$$w_{ji} := \frac{1}{m_j} \sum_{k=1}^{m_j} F_i(x_k) - \mathbb{E}_{q_i}[F_i];$$
 (13)

(b) Set

$$v_{ji} := \max(u_i, w_{ji} - 2^{-j-1}); \tag{14}$$

- (c) If  $\sum_{i} \min(2^{-j}, v_{ji} u_i) \cdot h_i > 2d'$ :
  - i. set  $v := v_j$
  - ii. terminate and output  $(j, v, \{F_i\})$
- 3. If the loop hasn't terminated for any j, output 'Fail' and restart the algorithm.

Figure 10: The Refined Hypothesis Select Algorithm.

 $\left(\frac{1}{m_j}\sum_{k=1}^{m_j}F_i(x_k)-\mathbb{E}_{q_i}[F_i]=v_i+2^{-j-1}\right)$ . Therefore, the event from the claim is equal to the event

$$(\mathbb{E}_p[F_i] - \mathbb{E}_{q_i}[F_i] \le v_i + 2^{-j-2}) \wedge \left( \frac{1}{m_j} \sum_{k=1}^{m_j} F_i(x_k) - \mathbb{E}_{q_i}[F_i] = v_i + 2^{-j-1} \right),$$

which is dominated by the event

$$\left| \mathbb{E}_p[F_i] - \frac{1}{m_j} \sum_{k=1}^{m_j} F_i(x_k) \right| \ge 2^{-j-2}.$$

By Chernoff bound, this probability is bounded by  $c_2\gamma/(\log d)$  for a small constant  $c_2$  (which depends on  $C_1$ ). By taking union bound on the different possible j's in the algorithm, we obtain an upper bound of  $\gamma$  on the failure probability.

Next – more importantly – we need to establish that the probability that the algorithm outputs 'Fail' is bounded by  $\gamma$  (First property of the algorithm).

Claim 15. The probability that the Refined Hypothesis Select algorithm outputs 'Fail' is  $< \gamma$ , where the randomness comes from the samples from P that it receives.

*Proof.* Our starting point is the fact that  $h \in \mathcal{H}_{\mathcal{P}_{\mathcal{Q}},d-d'}(u)$ , and therefore  $\min_{v \in \mathcal{P}_{\mathcal{Q}}} \{h \cdot v\} > h \cdot u + d - d'$ . Hence

$$\sum_{i \in [n]} h_i \cdot (\mathbb{E}_p[F_i] - \mathbb{E}_{q_i}[F_i] - u_i) > d - d'.$$

Partition the set of coordinates [n] as follows. Let

$$S_j := \{ i \in [n] : \mathbb{E}_p[F_i] - \mathbb{E}_{q_i}[F_i] - u_i \in (2^{-j}, 2^{-j+1}] \}$$
(16)

for  $j \in \{0, ..., 2 + \lceil \log(1 + 1/d) \rceil \}$ . Denote  $j_{max} := 2 + \lceil \log(1 + 1/d) \rceil \}$ . Note that the sets  $S_j$  are mutually disjoint. Some coordinates may belong to none of the sets, but only if  $\mathbb{E}_p[F_i] - \mathbb{E}_{q_i}[F_i] - u_i < 2^{-j_{max}}$ . We have

$$\sum_{j=0}^{j_{max}} \sum_{i \in S_j} h_i \cdot (\mathbb{E}_p[F_i] - \mathbb{E}_{q_i}[F_i] - u_i) >$$

$$\sum_{i \in [p]} h_i \cdot (\mathbb{E}_p[F_i] - \mathbb{E}_{q_i}[F_i] - u_i - 2^{-j_{max}}) > d - d' - 2^{-j_{max}} > \frac{d}{2}.$$

Therefore, there exists some j such that

$$\sum_{i \in S_i} h_i \cdot (\mathbb{E}_p[F_i] - \mathbb{E}_{q_i}[F_i] - u_i) > \frac{d}{2j_{max}}$$

Therefore, for any constant  $C_2 > 0$ , for a sufficiently large  $C_0$  the is a j such that

$$\sum_{i \in S_j} h_i \cdot (\mathbb{E}_p[F_i] - \mathbb{E}_{q_i}[F_i] - u_i) > C_2 \cdot d'. \tag{17}$$

Note that (17) and (16) implies

$$\sum_{i \in S_j} h_i > C_2 \cdot d' \cdot 2^j / 2. \tag{18}$$

We claim that for a sufficiently large constant  $C_2$ , the algorithm will terminate at step j with probability  $> 1 - \gamma$  (assuming it hasn't terminated earlier). Thus, the failure probability of the algorithm is bounded by  $\gamma$ .

For any given  $i \in S_j$ , we have by the Chernoff bound

$$\Pr[w_{ii} > (\mathbb{E}_n[F_i] - \mathbb{E}_{q_i}[F_i] - 2^{-j-2})] > 1 - \gamma/2,$$

and thus

$$\Pr[w_{ji} > u_i + 3 \cdot 2^{-j-2})] > 1 - \gamma/2.$$

Therefore

$$\mathbb{E}\left[\sum_{i \in S_j} h_i \cdot 1_{w_{ji} \le u_i + 3 \cdot 2^{-j-2}}\right] < \frac{\gamma}{2} \cdot \sum_{i \in S_j} h_i.$$

Therefore, by Markov inequality, with probability at least  $1 - \gamma$ ,

$$\sum_{i \in S_j} h_i \cdot 1_{w_{ji} \le u_i + 3 \cdot 2^{-j - 2}} < \frac{1}{2} \cdot \sum_{i \in S_j} h_i.$$

Hence, with probability at least  $1 - \gamma$ ,

$$\sum_{i \in S_i} h_i \cdot 1_{w_{ji} > u_i + 3 \cdot 2^{-j-2}} > \frac{1}{2} \cdot \sum_{i \in S_i} h_i. \tag{19}$$

We claim that assuming (19) holds, the algorithm will terminate at step j. We have

$$\sum_{i \in S_{j}: w_{ji} > u_{i} + 3 \cdot 2^{-j-2}} \min(2^{-j}, v_{i} - u_{i}) \cdot h_{i} = \sum_{i \in S_{j}: w_{ji} > u_{i} + 3 \cdot 2^{-j-2}} \min(2^{-j}, w_{ji} - 2^{-j-1}) \cdot h_{i} \geq \sum_{i \in S_{j}: w_{ji} > u_{i} + 3 \cdot 2^{-j-2}} 2^{-j-2} \cdot h_{i} \geq \frac{2^{-j-2}}{2} \cdot \sum_{i \in S_{j}} h_{i} \geq \sum_{i \in S_{j}: w_{ji} > u_{i} + 3 \cdot 2^{-j-2}} 2^{-j-2} \cdot h_{i} \geq \frac{2^{-j-2}}{2} \cdot C_{2} \cdot d' \cdot 2^{j} / 2 = \frac{C_{2} \cdot d'}{16} > 2d',$$

when  $C_2 > 32$  – guaranteeing that the algorithm terminates.

Claim 15 implies:

Claim 16. Assuming  $\gamma < d^3$ , the expected number of queries contributed by 'Fail's is an additive O(1).

*Proof.* The 'Fail' state is reached with probability  $< \gamma < d^3$ , while the number of queries of one run of the main loop is bounded by  $\tilde{O}(1/d^2)$ .

#### Hypothesis Selection in the Tiny Error Regime

- 1. Repeat the following until **Success** is reached:
  - (a) Run the refined Primal Game with  $\delta' = \varepsilon^2/n^3$  to obtain a distribution r;
  - (b) Use  $\tilde{O}(\log(1/\delta^2)/\varepsilon^2)$  fresh samples to verify that except with probability  $< \delta/2$ , for all calls of Hypothesis Selection Algorithm, whenever  $v_{ji} > u_i$ , we have  $\mathbb{E}_p[F_i] \mathbb{E}_{q_i}[F_i] > v_{ji}$ .
    - i. if verification passes, output **Success** and the distribution r;
    - ii. otherwise, restart the calculation.

Figure 11: The Tiny Error Case

### 6.3 Proof of Theorem 1

We are new ready to prove our main result, Theorem 1. We break the proof into two cases: the case when  $\delta$  is not too small:  $\delta \geq \varepsilon^2/n^3$ , and the case when  $\delta < \varepsilon^2/n^3$  is very small.

The case  $\delta \geq \varepsilon^2/n^3$ . In this case, we simply run the Refined Primal Game algorithm from Figure 9, where we set the error parameter  $\gamma$  in the Refined Hypothesis Select algorithm to  $\tilde{o}(\delta \varepsilon^2/\log n)$ .

**Correctness.** By Claim 13 applied to  $i^* := \operatorname{argmin}_i \mathsf{TV}(p, q_i)$ , we have, with probability  $> 1 - \delta$ , the output r satisfies

$$\mathsf{TV}(r, q_{i^*}) \le \mathsf{TV}(p, q_{i^*}) + \varepsilon/2,$$

therefore,

$$\mathsf{TV}(r,p) \le \mathsf{TV}(r,q_{i^*}) + \mathsf{TV}(p,q_{i^*}) < 2 \cdot \mathsf{TV}(p,q_{i^*}) + \varepsilon = 2 \cdot \min_i \mathsf{TV}(p,q_i) + \varepsilon. \tag{20}$$

**Sample complexity.** The conditions of Claim 12 are met with  $A = \tilde{O}(\log(1/\delta))$ . Therefore, by Claim 12, the total sample complexity in this case is bounded by

$$\tilde{O}\left(\frac{\log n \cdot \log(1/\delta)}{\varepsilon^2}\right). \tag{21}$$

The case  $\delta < \varepsilon^2/n^3$ . Consider the algorithm on Figure 11.

**Correctness.** The number of calls to Hypothesis Selection Algorithm is significantly smaller than  $o(1/\delta)$ . Therefore, by union bound, the probability of **Success** being returned despite  $\mathbb{E}_p[F_i] - \mathbb{E}_{q_i}[F_i] < v_{ji}$  holding at some point of the execution is  $o(\delta)$ . Assuming  $\mathbb{E}_p[F_i] - \mathbb{E}_{q_i}[F_i] > v_{ji}$  at all steps of the execution, the algorithm outputs a correct solution.

**Sample complexity.** In this case, we first run the previous case with  $\delta' = \varepsilon^2/n^3$ . As seen above, this step only requires

$$\tilde{O}\left(\frac{\log^2 n}{\varepsilon^2}\right)$$

samples. Moreover, as noted earlier, by union bound, with probability > 1 - 1/n, the event (15) from Claim 14 never happens throughout the execution of the algorithm. When (15) doesn't happen, we have

$$\mathbb{E}_p[F_i] - \mathbb{E}_{q_i}[F_i] > v_{ji} + 2^{-j-2} > v_{ji} + \tilde{\Omega}(\varepsilon),$$

and verification will pass with probability  $> 1 - \delta'$ . Therefore, the expected number of samples that will be needed until **Success** is reached is bounded by a

$$(1 + o(1)) \cdot (\text{number of samples used by one iteration}) = \tilde{O}\left(\frac{\log^2 n + \log(1/\delta)}{\varepsilon^2}\right).$$

# Acknowledgements

We thank Abbas Mehrabian and Hassan Zokaei Ashtiani for discussions regarding the implied improvement of our work to learning mixtures of gaussians. We also thank Naman Agarwal, Elad Hazan, Tomer Koren, and Karan Singh for fruitful discussions concerning the connections between the cutting-with-margin game and online optimization.

## References

- [ABDH<sup>+</sup>20] Hassan Ashtiani, Shai Ben-David, Nicholas J. A. Harvey, Christopher Liaw, Abbas Mehrabian, and Yaniv Plan. Near-optimal sample complexity bounds for robust learning of gaussian mixtures via compression schemes. *J. ACM*, 67(6), October 2020.
- [ABM18] Hassan Ashtiani, Shai Ben-David, and Abbas Mehrabian. Sample-efficient learning of mixtures. In *Conference on Artificial Intelligence (AAAI)*, pages 2679–2686, 2018.
- [AFJ<sup>+</sup>18] Jayadev Acharya, Moein Falahatgar, Ashkan Jafarpour, Alon Orlitsky, and Ananda Theertha Suresh. Maximum selection and sorting with adversarial comparators. J. Mach. Learn. Res., 19:59:1–59:31, 2018.
- [AHK12] Sanjeev Arora, Elad Hazan, and Satyen Kale. The multiplicative weights update method: a meta-algorithm and applications. *Theory of Computing*, 8(6):121–164, 2012.
- [AJOS14] Jayadev Acharya, Ashkan Jafarpour, Alon Orlitsky, and Ananda Theertha Suresh. Sorting with adversarial comparators and application to density estimation. In *International Symposium on Information Theory (ISIT)*, pages 1682–1686. IEEE, 2014.
- [AW01] Katy S. Azoury and M. K. Warmuth. Relative loss bounds for on-line density estimation with the exponential family of distributions. *Mach. Learn.*, 43(3):211–246, June 2001.

- [BKM19] Olivier Bousquet, Daniel Kane, and Shay Moran. In Conference on Learning Theory (COLT), volume 99, pages 318–341, 2019.
- [BKSW21] Mark Bun, Gautam Kamath, Thomas Steinke, and Zhiwei Steven Wu. Private hypothesis selection. *IEEE Trans. Inf. Theory*, 67(3):1981–2000, 2021.
- [Bub15] Sébastien Bubeck. Convex optimization: Algorithms and complexity. Found. Trends Mach. Learn., 8(3–4):231–357, November 2015.
- [CDSS14] Siu-on Chan, Ilias Diakonikolas, Rocco A. Servedio, and Xiaorui Sun. Near-optimal density estimation in near-linear time using variable-width histograms. In *Advances in Neural Information Processing Systems (NIPS)*, pages 1844–1852, 2014.
- [DDS15] Constantinos Daskalakis, Ilias Diakonikolas, and Rocco A. Servedio. Learning poisson binomial distributions. *Algorithmica*, 72(1):316–357, 2015.
- [DFH<sup>+</sup>15] Cynthia Dwork, Vitaly Feldman, Moritz Hardt, Toniann Pitassi, Omer Reingold, and Aaron Leon Roth. Preserving statistical validity in adaptive data analysis. In Symposium on Theory of Computing (STOC), pages 117–126. ACM, 2015.
- [Dia16] Ilias Diakonikolas. Learning structured distributions. In *Handbook of Big Data*, pages 267–283. Chapman and Hall/CRC, 2016.
- [DK14] Constantinos Daskalakis and Gautam Kamath. Faster and sample near-optimal algorithms for proper learning mixtures of gaussians. In *Conference on Learning Theory (COLT)*, volume 35 of *JMLR Workshop and Conference Proceedings*, pages 1183–1213, 2014.
- [DKK<sup>+</sup>19] Ilias Diakonikolas, Gautam Kamath, Daniel Kane, Jerry Li, Ankur Moitra, and Alistair Stewart. Robust estimators in high-dimensions without the computational intractability. SIAM J. Comput., 48(2):742–864, 2019.
- [DKS17] Ilias Diakonikolas, Daniel M. Kane, and Alistair Stewart. Statistical query lower bounds for robust estimation of high-dimensional gaussians and gaussian mixtures. In Foundations of Computer Science (FOCS), pages 73–84, 2017.
- [DKS18] Ilias Diakonikolas, Daniel M. Kane, and Alistair Stewart. List-decodable robust mean estimation and learning mixtures of spherical gaussians. In *Symposium on Theory of Computing (STOC)*, pages 1047–1060. ACM, 2018.
- [DL96] Luc Devroye and Gábor Lugosi. A universally acceptable smoothing factor for kernel density estimates. *The Annals of Statistics*, pages 2499–2512, 1996.
- [DL97] Luc Devroye and Gábor Lugosi. Nonasymptotic universal smoothing factors, kernel complexity and yatracos classes. *The Annals of Statistics*, 25(6):2626–2637, 1997.
- [DL01] L. Devroye and G. Lugosi. Combinatorial methods in density estimation. Springer, 2001.
- [DL04] Luc Devroye and Gábor Lugosi. Bin width selection in multivariate histograms by the combinatorial method. *Test*, 13(1):129–145, Jun 2004.
- [DLS18] Ilias Diakonikolas, Jerry Li, and Ludwig Schmidt. Fast and sample near-optimal algorithms for learning multidimensional histograms. In Sébastien Bubeck, Vianney

- Perchet, and Philippe Rigollet, editors, *Proceedings of the 31st Conference On Learning Theory*, volume 75 of *Proceedings of Machine Learning Research*, pages 819–842. PMLR, 06–09 Jul 2018.
- [DS14] Amit Daniely and Shai Shalev-Shwartz. Optimal learners for multiclass problems. In Maria-Florina Balcan, Vitaly Feldman, and Csaba Szepesvári, editors, *Conference on Learning Theory (COLT)*, volume 35, pages 287–316, 2014.
- [GKK+20] Sivakanth Gopi, Gautam Kamath, Janardhan Kulkarni, Aleksandar Nikolov, Zhiwei Steven Wu, and Huanyu Zhang. Locally private hypothesis selection. In Conference on Learning Theory (COLT), volume 125 of Proceedings of Machine Learning Research, pages 1785–1816, 2020.
- [Grü60] B. Grünbaum. Partitions of mass-distributions and of convex bodies by hyperplanes. Pacific J. Math., 10(4):1257–1261, 1960.
- [JHW18] J. Jiao, Y. Han, and T. Weissman. Minimax estimation of the  $l_1$  distance. *IEEE Transactions on Information Theory*, 64(10):6672–6706, Oct 2018.
- [KMV12] Adam Tauman Kalai, Ankur Moitra, and Gregory Valiant. Disentangling gaussians. Commun. ACM, 55(2):113–120, 2012.
- [KSS18] Pravesh K. Kothari, Jacob Steinhardt, and David Steurer. Robust moment estimation and improved clustering via sum of squares. In *Symposium on Theory of Computing* (STOC), pages 1035–1046, 2018.
- [LN96] Gábor Lugosi and Andrew Nobel. Consistency of data-driven histogram methods for density estimation and classification. *Ann. Statist.*, 24(2):687–706, 04 1996.
- [MS08] Satyaki Mahalanabis and Daniel Stefankovic. Density estimation in linear time. In Rocco A. Servedio and Tong Zhang, editors, Conference on Learning Theory (COLT), pages 503–512, 2008.
- [Pea95] K. Pearson. Contributions to the mathematical theory of evolution. ii. skew variation in homogeneous material. *Philosophical Trans. of the Royal Society of London*, 186:343–414, 1895.
- [PW15] Yury Polyanskiy and Yihong Wu. Lecture notes on information theory, 2015.
- [SF12] Robert E Schapire and Yoav Freund. Boosting: Foundations and algorithms. MIT press, 2012.
- [SOAJ14] Ananda Theertha Suresh, Alon Orlitsky, Jayadev Acharya, and Ashkan Jafarpour. Near-optimal-sample estimators for spherical gaussian mixtures. In *Advances in Neural Information Processing Systems (NIPS)*, pages 1395–1403, 2014.
- [vN28] J. von Neumann. Zur theorie der gesellschaftsspiele. *Mathematische Annalen*, 100:295–320, 1928.
- [Yat85] Yannis G. Yatracos. Ann. Statist., 13(2):768–774, 06 1985.