

# Agnostic Learning of Monomials by Halfspaces is Hard\*

Vitaly Feldman<sup>†</sup> Venkatesan Guruswami<sup>‡</sup> Prasad Raghavendra<sup>§</sup> Yi Wu<sup>¶</sup>

December 1, 2010

## Abstract

We prove the following strong hardness result for learning: Given a distribution of labeled examples from the hypercube such that there exists a *monomial* consistent with  $(1 - \epsilon)$  of the examples, it is NP-hard to find a *halfspace* that is correct on  $(1/2 + \epsilon)$  of the examples, for arbitrary constants  $\epsilon > 0$ . In learning theory terms, weak agnostic learning of monomials is hard, even if one is allowed to output a hypothesis from the much bigger concept class of halfspaces. This hardness result subsumes a long line of previous results, including two recent hardness results for the proper learning of monomials and halfspaces. As an immediate corollary of our result we show that weak agnostic learning of *decision lists* is NP-hard.

Our techniques are quite different from previous hardness proofs for learning. We define distributions on positive and negative examples for monomials whose first few moments match. We use the *invariance principle* to argue that *regular* halfspaces (all of whose coefficients have small absolute value relative to the total  $\ell_2$  norm) cannot distinguish between distributions whose first few moments match. For highly non-regular subspaces, we use a structural lemma from recent work on fooling halfspaces to argue that they are “junta-like” and one can zero out all but the top few coefficients without affecting the performance of the halfspace. The top few coefficients form the natural list decoding of a halfspace in the context of dictatorship tests/Label Cover reductions.

We note that unlike previous invariance principle based proofs which are only known to give Unique-Games hardness, we are able to reduce from a version of Label Cover problem that is known to be NP-hard. This has inspired follow-up work on bypassing the Unique Games conjecture in some optimal geometric inapproximability results.

---

\*An extended abstract appeared in the Proceedings of the 50th IEEE Symposium on Foundations of Computer Science, 2009.

<sup>†</sup>IBM Almaden Research Center, San Jose, CA. [vitaly@post.harvard.edu](mailto:vitaly@post.harvard.edu).

<sup>‡</sup>Computer Science Department, Carnegie Mellon University, Pittsburgh, PA. [guruswami@cmu.edu](mailto:guruswami@cmu.edu).

<sup>§</sup>College of Computing, Georgia Institute of Technology, Atlanta, GA. [praghava@cc.gatech.edu](mailto:praghava@cc.gatech.edu). Some of this work was done when visiting Carnegie Mellon University.

<sup>¶</sup>IBM Almaden Research Center, San Jose, CA. [wuyi@us.ibm.com](mailto:wuyi@us.ibm.com). Most of this work was done when the author was at Carnegie Mellon University.

# 1 Introduction

Boolean conjunctions (or *monomials*), decision lists, and halfspaces are among the most basic concept classes in learning theory. They are all long-known to be efficiently PAC learnable, when the given examples are guaranteed to be consistent with a function from any of these concept classes [42, 7, 39]. However, in practice data is often noisy or too complex to be consistently explained by a simple concept. A common practical approach to such problems is to find a predictor in a certain space of hypotheses that best fits the given examples. A general model for learning that addresses this scenario is the *agnostic* learning model [21, 26]. An *agnostic* learning algorithm for a class of functions  $\mathcal{C}$  using a hypothesis space  $\mathcal{H}$  is required to perform the following task: Given examples drawn from some unknown distribution, the algorithm must find a hypothesis in  $\mathcal{H}$  that classifies the examples nearly as well as is possible by a hypothesis from  $\mathcal{C}$ . The algorithm is said to be a *proper* learning algorithm if  $\mathcal{C} = \mathcal{H}$ .

In this work we address the complexity of agnostic learning of monomials by algorithms that output a halfspace as a hypothesis. Learning methods that output a halfspace as a hypothesis such as Perceptron [40], Winnow [34], Support Vector Machines [43] as well as most boosting algorithms are well-studied in theory and widely used in practical prediction systems. These classifiers are often applied to labeled data sets which are not linearly separable. Hence it is of great interest to determine the classes of problems that can be solved by such methods in the agnostic setting. In this work we demonstrate a strong negative result on agnostic learning by halfspaces. We prove that non-trivial agnostic learning of even the relatively simple class of monomials by halfspaces is an NP-hard problem.

**Theorem 1.1.** *For any constant  $\epsilon > 0$ , it is NP-hard to find a halfspace that correctly labels  $(1/2 + \epsilon)$ -fraction of given examples over  $\{0, 1\}^n$  even when there exists a monomial that agrees with a  $(1 - \epsilon)$ -fraction of the examples.*

Note that this hardness result is essentially optimal since it is trivial to find a hypothesis with agreement rate  $1/2$  — output either the function that is always 0 or the function that is always 1. Also note that Theorem 1.1 measures agreement of a halfspace and a monomial with the given set of examples rather than the probability of agreement of  $h$  with an example drawn randomly from an unknown distribution. Uniform convergence results based on the VC dimension imply that these settings are essentially equivalent (see for example [21, 26]).

The class of monomials is a subset of the class of decision lists which in turn is a subset of the class of halfspaces. Therefore our result immediately implies an optimal hardness result for proper agnostic learning of decision lists.

## Previous work

Before describing the details of the prior body of work on hardness results for learning, we note that our result *subsumes all these results* with just one exception (the hardness of learning monomials by  $t$ -CNFs [32]). This is because we obtain the optimal inapproximability factor *and* allow learning of monomials by the much richer class of halfspaces.

The results of the paper are noteworthy in the broader context of hardness of approximation. Previously, hardness proofs based on the invariance principle were only known to give Unique-Games

hardness. In this work, we are able to harness invariance principles to show NP-hardness result by working with a version of Label Cover whose projection functions are only required to be *unique-on-average*. This could be one potential approach to revisit the many strong inapproximability results conditioned on the Unique Games conjecture (UGC), with an eye towards bypassing the UGC assumption. Such a goal was achieved for some geometric problems recently [?]; see Section 2.3.

Agnostic learning of monomials, decision lists and halfspaces has been studied in a number of previous works. Proper agnostic learning of a class of functions  $\mathcal{C}$  is equivalent to the ability to come up with a function in  $\mathcal{C}$  which has the optimal agreement rate with the given set of examples and is also referred to as the *Maximum Agreement* problem for a class of function  $\mathcal{C}$ .

The Maximum Agreement problem for halfspaces is equivalent to the so-called Hemisphere problem and is long known to be NP-complete [23, 17]. Amaldi and Kann [1] showed that Maximum Agreement for halfspaces is NP-hard to approximate within  $\frac{261}{262}$  factor. This was later improved by Ben-David *et al.* [5], and Bshouty and Burroughs [9] to approximation factors  $\frac{415}{418}$ , and  $\frac{84}{85}$ , respectively. An optimal inapproximability result was established independently by Guruswami and Raghavendra [20] and Feldman *et al.* [15] showing NP-hardness of approximating the Maximum Agreement problem for halfspaces within  $(1/2 + \epsilon)$  for every constant  $\epsilon > 0$ . The reduction in [15] requires examples with real-valued coordinates, whereas the proof in [20] also works for examples drawn from the Boolean hypercube.

The Maximum Agreement problem for monotone monomials was shown to be NP-hard by Angluin and Laird [2], and NP-hardness for general monomials was shown by Kearns and Li [27]. The hardness of approximating the maximum agreement within  $\frac{767}{770}$  was shown by Ben-David *et al.* [5]. The factor was subsequently improved to  $58/59$  by Bshouty and Burroughs [9]. Finally, Feldman *et al.* [14, 15] showed a tight inapproximability result, namely that it is NP-hard to distinguish between the instances where  $(1 - \epsilon)$ -fraction of the labeled examples are consistent with some monomial and instances where every monomial is consistent with at most  $(1/2 + \epsilon)$ -fraction of the examples. Recently, Khot and Saket [32] proved a similar hardness result even when a  $t$ -CNF is allowed as output hypothesis for an arbitrary constant  $t$  (a  $t$ -CNF is the conjunction of several clauses, each of which has at most  $t$  literals; a monomial is thus a 1-CNF).

For the concept class of decisions lists, APX-hardness (or hardness to approximate within some constant factor) of the Maximum Agreement problem was shown by Bshouty and Burroughs [9]. As mentioned above, our result subsumes all these results with the exception of [32].

A number of hardness of approximation results are also known for the complementary problem of minimizing disagreement for each of the above concept classes [26, 22, 3, 8, 14, 15]. Another well-known evidence of the hardness of agnostic learning of monomials is that even a non-proper agnostic learning of monomials would give an algorithm for learning DNF — a major open problem in learning theory [33]. Further, Kalai *et al.* proved that even agnostic learning of halfspaces with respect to the uniform distribution implies learning of parities with random classification noise — a long-standing open problem in learning theory and coding [24].

Monomials, decision lists and halfspaces are known to be efficiently learnable in the presence of more benign *random* classification noise [2, 25, 28, 10, 6, 12]. Simple online algorithms like Perceptron and Winnow learn halfspaces when the examples can be separated with a significant *margin* (as is the case if the examples are consistent with a monomial) and are known to be robust

to a very mild amount of adversarial noise [16, 4, 18]. Our result implies that these positive results will not hold when the adversarial noise rate is  $\epsilon$  for any constant  $\epsilon > 0$ .

Kalai *et al.* gave the first non-trivial algorithm for agnostic learning monomials in time  $2^{\tilde{O}(\sqrt{n})}$  [24]. They also gave a breakthrough result for agnostic learning of halfspaces with respect to the uniform distribution on the hypercube up to any constant accuracy (and analogous results for a number of other settings). Their algorithms output linear thresholds of parities as hypotheses. In contrast, our hardness result is for algorithms that output a halfspace (which is a linear threshold of single variables).

**Organization of the paper:** We sketch the idea of our proof in Section 2. We define some probability and analytical tools in Section 3. In Section 4 we define the *dictatorship test*, which is an important gadget for the hardness reduction. For the purpose of illustration, we also show why this dictatorship test already suffices to prove Theorem 1.1 assuming the Unique Games Conjecture [29]. In Section 5, we describe a reduction from a variant of the LABEL COVER problem to prove Theorem 1.1 under the assumption that  $P \neq NP$ .

**Notation:** We use 0 to encode “False” and 1 to encode “True”. We denote  $\text{pos}(t) : \mathbb{R} \rightarrow \{0, 1\}$  as the indicator function of whether  $t \geq 0$ ; i.e.,  $\text{pos}(t) = 1$  when  $t \geq 0$  and  $\text{pos}(t) = 0$  when  $t < 0$ .

For  $\mathbf{x} = (x_1, x_2, \dots, x_n) \in \{0, 1\}^n$ ,  $\mathbf{w} \in \mathbb{R}^n$ , and  $\theta \in \mathbb{R}$ , a halfspace  $h(\mathbf{x})$  is a Boolean function of the form  $\text{pos}(\mathbf{w} \cdot \mathbf{x} - \theta)$ ; a monomial (conjunction) is a function of the form  $\bigwedge_{i \in S} s_i$ , where  $S \subseteq [n]$  and  $s_i$  is the literal of  $x_i$  which can represent either  $x_i$  or  $\neg x_i$ ; a disjunction is a function of the form  $\bigvee_{i \in S} s_i$ . One special case of monomials is the function  $f(x) = x_i$  for some  $i \in [n]$ , also referred to as the  $i$ -th *dictator* function.

## 2 Proof Overview

We prove Theorem 1.1 by exhibiting a reduction from the  $k$ -LABEL COVER problem, which is a particular variant of the LABEL COVER problem. The  $k$ -LABEL COVER problem is defined as follows:

**Definition 2.1.** For positive integer  $M, N$  that  $M \geq N$  and  $k \geq 2$ , an instance of  $k$ -LABEL COVER  $\mathcal{L}(G(V, E), M, N, \{\pi^{v,e} | e \in E, v \in e\})$  consists of a  $k$ -uniform connected (multi-)hypergraph  $G(V, E)$  with vertex set  $V$  and an edge multiset  $E$ ; a set of functions  $\{\pi^{v_i,e}\}_{i=1}^k$ . Every hyperedge  $e = (v_1, \dots, v_k)$  is associated with a  $k$ -tuple of projection functions  $\{\pi^{v_i,e}\}_{i=1}^k$  where  $\pi^{v_i,e} : [M] \rightarrow [N]$ .

A vertex labeling  $\Lambda$  is an assignment of labels to vertices  $\Lambda : V \rightarrow [M]$ . A labeling  $\Lambda$  is said to strongly satisfy an edge  $e$  if  $\pi^{v_i,e}(\Lambda(v_i)) = \pi^{v_j,e}(\Lambda(v_j))$  for every  $v_i, v_j \in e$ . A labeling  $\Lambda$  weakly satisfies edge  $e$  if  $\pi^{v_i,e}(\Lambda(v_i)) = \pi^{v_j,e}(\Lambda(v_j))$  for some  $v_i, v_j \in e, v_i \neq v_j$ .

The goal in LABEL COVER is to find a vertex labeling that satisfies as many edges (projection constraints) as possible.

## 2.1 Hardness assuming the Unique Games conjecture

For the sake of clarity, we first sketch the proof of Theorem 1.1 with a reduction from the  $k$ -UNIQUE LABEL COVER problem which is a special case of  $k$ -LABEL COVER where  $M = N$  and all the projection functions  $\{\pi^{v,e} | v \in e, e \in E\}$  are bijections. The following inapproximability result [?] for  $k$ -UNIQUE LABEL COVER is equivalent to the Unique Games Conjecture of Khot [29].

**Conjecture 2.2.** *For every constant  $\eta > 0$  and a positive integer  $k$ , there exists an integer  $R_0$  such that for all positive integers  $R > R_0$ , given an instance  $\mathcal{L}(G(V, E), R, R, \{\pi^{v,e} | e \in E, v \in e\})$  it is NP-hard to distinguish between,*

- *strongly satisfiable instances: there exists a labeling  $\Lambda : V \rightarrow [R]$  that strongly satisfies  $1 - k\eta$  fraction of the edges  $E$ .*
- *almost unsatisfiable instances: there is no labeling that weakly satisfies  $\frac{2k^2}{R^{\eta/4}}$  fraction of the edges.*

Given an instance  $\mathcal{L}$  of  $k$ -UNIQUE LABEL COVER, we will produce a distribution  $\mathcal{D}$  over labeled examples such that the following holds: if  $\mathcal{L}$  is a strongly satisfiable instance, then there is a disjunction that agrees with the label on a randomly chosen example with probability at least  $1 - \epsilon$ , while if  $\mathcal{L}$  is an almost unsatisfiable instance then no halfspace agrees with the label on a random example from  $\mathcal{D}$  with probability more than  $\frac{1}{2} + \epsilon$ . Clearly, such a reduction implies Theorem 1.1 assuming the Unique Games Conjecture but with disjunctions in place of conjunctions. De Morgan's law and the fact that a negation of a halfspace is a halfspace then imply that the statement is also true for monomials (we use disjunctions only for convenience).

Let  $\mathcal{L}$  be an instance of  $k$ -UNIQUE LABEL COVER on hypergraph  $G = (V, E)$  and a set of labels  $[R]$ . The examples we generate will have  $|V| \times R$  coordinates, i.e., belong to  $\{0, 1\}^{|V| \times R}$ . These coordinates are to be thought of as one block of  $R$  coordinates for every vertex  $v \in V$ . We will index the coordinates of  $\mathbf{x} \in \{0, 1\}^{|V| \times R}$  as  $\mathbf{x} = (x_v^{(r)})_{v \in V, r \in [R]}$ .

For every labeling  $\Lambda : V \rightarrow [R]$  of the instance, there is a corresponding disjunction over  $\{0, 1\}^{|V| \times R}$  given by,

$$h(\mathbf{x}) = \bigvee_v x_v^{(\Lambda(v))}.$$

Thus, using a label  $r$  for a vertex  $v$  is encoded as including the literal  $x_v^{(r)}$  in the disjunction. Notice that an arbitrary halfspace over  $\{0, 1\}^{|V| \times R}$  need not correspond to any labeling at all. The idea would be to construct a distribution on examples which ensures that any halfspace agreeing with at least  $\frac{1}{2} + \epsilon$  fraction of random examples somehow corresponds to a labeling of  $\Lambda$  weakly satisfying a constant fraction of the edges in  $\mathcal{L}$ .

Fix an edge  $e = (v_1, \dots, v_k)$ . For the sake of exposition, let us assume  $\pi^{v_i, e}$  is the identity permutation for every  $i \in [k]$ . The general case is not anymore complicated.

For the edge  $e$ , we will construct a distribution on examples  $\mathcal{D}_e$  with the following properties:

- All coordinates  $x_v^{(r)}$  for a vertex  $v \notin e$  are fixed to be zero. Restricted to these examples, the halfspace  $h$  can be written as  $h(\mathbf{x}) = \text{pos}(\sum_{i \in [k]} \langle \mathbf{w}_{v_i}, \mathbf{x}_{v_i} \rangle - \theta)$ .

- For any label  $r \in [R]$ , the labeling  $\Lambda(v_1) = \dots = \Lambda(v_k) = r$  *strongly* satisfies the edge  $e$ . Hence, the corresponding disjunction  $\bigvee_{i \in [k]} x_{v_i}^{(r)}$  needs to have agreement  $\geq 1 - \epsilon$  with the examples from  $\mathcal{D}_e$ .
- There exists a decoding procedure that given a halfspace  $h$  outputs a labeling  $\Lambda_h$  for  $\mathcal{L}$  such that, if  $h$  has agreement  $\geq \frac{1}{2} + \epsilon$  with the examples from  $\mathcal{D}_e$ , then  $\Lambda_h$  *weakly* satisfies the edge  $e$  with non-negligible probability.

For conceptual clarity, let us rephrase the above requirement as a testing problem. Given a halfspace  $h$ , consider a randomized procedure that samples an example  $(\mathbf{x}, b)$  from the distribution  $\mathcal{D}_e$ , and accepts if  $h(\mathbf{x}) = b$ . This amounts to a test that checks if the function  $h$  corresponds to a consistent labeling. Further, let us suppose the halfspace  $h$  is given by  $h(\mathbf{x}) = \text{pos}(\sum_{v \in V} \langle \mathbf{w}_v, \mathbf{x}_v \rangle - \theta)$ . Define the linear function  $f_v : \{0, 1\}^R \rightarrow \mathbb{R}$  as  $f_v(\mathbf{x}_v) = \langle \mathbf{w}_v, \mathbf{x}_v \rangle$ . Then, we have  $h(\mathbf{x}) = \text{pos}(\sum_{v \in V} f_v(\mathbf{x}_v) - \theta)$ .

For a halfspace  $h$  corresponding to a labeling  $\Lambda$ , we will have  $f_v(\mathbf{x}_v) = x_v^{(\Lambda(v))}$  – a dictator function. Thus, in the intended solution every linear function  $f_v$  associated with the halfspace  $h$  is a dictator function.

Now, let us again restate the above testing problem in terms of these linear functions. For succinctness, we write  $f_i$  for the linear function  $f_{v_i}$ . We need a randomized procedure that does the following:

Given  $k$  linear functions  $f_1, \dots, f_k : \{0, 1\}^R \rightarrow \mathbb{R}$ , queries the functions at one point each (say  $\mathbf{x}_1, \dots, \mathbf{x}_k$  respectively), and accepts if  $\text{pos}(\sum_{i=1}^k f_i(\mathbf{x}_i) - \theta) = b$ .

The procedure must satisfy,

- (Completeness) If each of the linear functions  $f_i$  is the  $r$ 'th dictator function for some  $r \in [R]$ , then the test accepts with probability  $1 - \epsilon$ .
- (Soundness) If the test accepts with probability  $\frac{1}{2} + \epsilon$ , then at least *two* of the linear functions are *close* to the same dictator function.

A testing problem of the above nature is referred to as a *Dictatorship Testing* and is a recurring theme in hardness of approximation.

Notice that the notion of a linear function being *close* to a dictator function is not formally defined yet. In most applications, a function is said to be close to a dictator if it has *influential* coordinates. It is easy to see that this notion is not sufficient by itself here. For example, in the linear function  $\text{pos}(10^{100}x_1 + x_2 - 0.5)$ , although the coordinate  $x_2$  has little influence on the linear function, it has significant influence on the halfspace.

We resolve this problem by using the notion of *critical index* (Definition 3.1) that was introduced in [41] and has found numerous applications in the analysis of halfspaces [35, 38, 13]. Roughly speaking, given a linear function  $f$ , the idea is to recursively delete its influential coordinates until there are none left. The total number of coordinates so deleted is referred to as the critical index of  $f$ . Let  $c_\tau(\mathbf{w}_i)$  denote the critical index of  $\mathbf{w}_i$ , and let  $C_\tau(\mathbf{w}_i)$  denote the set of  $c_\tau(\mathbf{w}_i)$  largest coordinates of  $\mathbf{w}_i$ . The linear function  $l$  is said to be *close* to the  $i$ 'th dictator function for every  $i$

in  $C_\tau(\mathbf{w}_i)$ . A function is *far* from every dictator if it has critical index 0 – no influential coordinate to delete.

An important issue is that the critical index of a linear function can be much larger than the number of influential coordinates and cannot be appropriately bounded. In other words, a linear function can be close to a large number of dictator functions, as per the definition above. To counter this, we employ a structural lemma about halfspaces that was used in the recent work on fooling halfspaces with limited independence [13]. Using this lemma, we are able to prove that if the critical index is large, then one can in fact zero out the coordinates of  $\mathbf{w}_i$  outside the  $t$  largest coordinates for some large enough  $t$ , and the agreement of the halfspace  $h$  only changes by a negligible amount! Thus, we first carry out the zeroing operation for all linear functions with large critical index.

We now describe the above construction and analysis of the dictatorship test in some more detail. It is convenient to think of the  $k$  queries  $\mathbf{x}_1, \dots, \mathbf{x}_k$  as the rows of a  $k \times R$  matrix with  $\{0, 1\}$  entries. Henceforth, we will refer to matrices  $\{0, 1\}^{k \times R}$  and their rows and columns.

We construct two distributions  $\mathcal{D}_0, \mathcal{D}_1$  on  $\{0, 1\}^k$  such that for  $s \in \{0, 1\}$ , we have  $\Pr_{\mathbf{x} \in \mathcal{D}_s} [\bigvee_{i=1}^k x_i = s] \geq 1 - \epsilon/2$  for  $\epsilon = o_k(1)$  (this will ensure the completeness of the reduction, i.e., certain disjunctions pass with high probability). Further, the distributions  $\mathcal{D}_0, \mathcal{D}_1$  will be carefully chosen to have matching first four moments. This will be used in the soundness analysis where we will use an *invariance principle* to infer structural properties of halfspaces that pass the test with probability noticeably greater than  $1/2$ .

We define the distribution  $\tilde{\mathcal{D}}_s^R$  on matrices  $\{0, 1\}^{k \times R}$  by sampling  $R$  columns independently according to  $\mathcal{D}_s$ , and then perturbing each bit with a small probability  $\epsilon/2$ . We define the following test (or equivalently, distribution on examples): given a halfspace  $h$  on  $\{0, 1\}^{k \times R}$ , with probability  $1/2$  we check  $h(\mathbf{x}) = 0$  for a sample  $\mathbf{x} \in \tilde{\mathcal{D}}_0^R$ , and with probability  $1/2$  we check  $h(\mathbf{x}) = 1$  for a sample  $\mathbf{x} \in \tilde{\mathcal{D}}_1^R$ .

**Completeness:** By construction, each of the  $R$  disjunctions  $\text{OR}_j(\mathbf{x}) = \bigvee_{i=1}^k x_i^{(j)}$  passes the test with probability at least  $1 - \epsilon$  (here  $x_i^{(j)}$  denotes the entry in the  $i$ 'th row and  $j$ 'th column of  $\mathbf{x}$ ).

**Soundness:** For the soundness analysis, suppose  $h(\mathbf{x}) = \text{pos}(\langle \mathbf{w}, \mathbf{x} \rangle - \theta)$  is a halfspace that passes the test with probability at least  $1/2 + \epsilon$ . The halfspace  $h$  can be written in two ways by expanding the inner product  $\langle \mathbf{w}, \mathbf{x} \rangle$  along rows and columns, i.e.,  $h(\mathbf{x}) = \text{pos}(\sum_{i=1}^k \langle \mathbf{w}_i, \mathbf{x}_i \rangle - \theta) = \text{pos}(\sum_{i=1}^R \langle \mathbf{w}^{(i)}, \mathbf{x}^{(i)} \rangle - \theta)$ . Let us denote  $f_i(\mathbf{x}) = \langle \mathbf{w}_i, \mathbf{x}_i \rangle$ .

First, let us see why the linear functions  $\langle \mathbf{w}_i, \mathbf{x}_i \rangle$  must be close to *some* dictator. Note that we need to show that two of the linear functions are close to the *same* dictator.

Suppose each of the linear functions  $f_i$  is not *close* to any dictator. In other words, for each  $i$ , no single coordinate of the vector  $\mathbf{w}_i$  is too large (contains more than  $\tau$ -fraction of the  $\ell_2$  mass  $\|\mathbf{w}_i\|_2$  of vector  $\mathbf{w}_i$ ). Clearly, this implies that no single column of the matrix  $\mathbf{w}$  is too *large*.

Recall that the halfspace is given by  $h(\mathbf{x}) = \text{pos}(\sum_{j \in [R]} \langle \mathbf{w}^{(j)}, \mathbf{x}^{(j)} \rangle - \theta)$ . Here  $l(\mathbf{x}) = \sum_{j \in [R]} \langle \mathbf{w}^{(j)}, \mathbf{x}^{(j)} \rangle - \theta$  is a degree 1 polynomial into which we are substituting values from two product distributions  $\mathcal{D}_0^R$  and  $\mathcal{D}_1^R$ . Further, the distributions  $\mathcal{D}_0$  and  $\mathcal{D}_1$  have matching moments up to order 4 by design. Using the invariance principle, the distribution of  $l(\mathbf{x})$  is roughly the same, whether  $\mathbf{x}$  is from  $\mathcal{D}_0^R$  or  $\mathcal{D}_1^R$ . Thus, by the invariance principle, the halfspace  $h$  is unable to distinguish between the distributions  $\mathcal{D}_0^R$  and  $\mathcal{D}_1^R$  with a noticeable advantage.

Further, suppose no two linear functions  $f_i$  are *close* to the same dictator, i.e.,  $C_\tau(\mathbf{w}_i) \cap C_\tau(\mathbf{w}_j) = \emptyset$ . In this case, we condition on the values of  $x_i^{(j)}$  for  $j \in C_\tau(\mathbf{w}_i)$ . Since  $C_\tau(\mathbf{w}_i) \cap C_\tau(\mathbf{w}_j) = \emptyset$ , this conditions at most *one* value in each column. Therefore, the conditional distribution on each column in cases  $\mathcal{D}_0$  and  $\mathcal{D}_1$  still have matching first three moments. We thus apply the invariance principle using the fact that after deleting the coordinates in  $C_\tau(\mathbf{w}_i)$ , all the remaining coefficients of the weight vector  $\mathbf{w}$  are small (by definition of critical index). This implies that  $C_\tau(\mathbf{w}_i) \cap C_\tau(\mathbf{w}_j) \neq \emptyset$  for some two rows  $i, j$  and finishes the proof of the soundness claim.

The above consistency-enforcing test almost immediately yields the Unique Games hardness of weak learning disjunctions by halfspaces via standard methods.

## 2.2 Extending to NP-hardness

To prove NP-hardness as opposed to hardness assuming the Unique Games conjecture, we reduce a version of Label Cover to our problem. This requires a more complicated consistency check, and we have to overcome several additional technical obstacles in the proof.

The main obstacle encountered in transferring the dictatorship test to a Label Cover-based hardness is one that commonly arises for several other problems. Specifically, the projection constraint on an edge  $e = (u, v)$  maps a large set of labels  $\mathcal{R} = \{r_1, \dots, r_d\}$  corresponding to a vertex  $u$  to a single label  $r$  for the vertex  $v$ . While composing the Label Cover constraint  $(u, v)$  with the dictatorship test, all labels in  $\mathcal{R}$  have to be necessarily *equivalent*. In several settings including this work, this requires the coordinates corresponding to labels in  $\mathcal{R}$  to be mostly identical! However, on making the coordinates corresponding to  $\mathcal{R}$  identical, the prover corresponding to  $u$  can determine the identity of edge  $(u, v)$ , thus completely destroying the soundness of the composition. In fact, the natural extension of the Unique Games-based reduction for MAXCUT [31] to a corresponding Label Cover hardness fails primarily for this reason.

Unlike MAXCUT or other Unique Games-based reductions, in our case, the soundness of the dictatorship test is required to hold against a specific class of functions, i.e, halfspaces. Harnessing this fact, we execute the reduction starting from a Label Cover instance whose projections are *unique on average*. More precisely, a *smooth* Label Cover (introduced in [30]) is one in which for every vertex  $u$ , and a pair of labels  $r, r'$ , the labels  $\{r, r'\}$  project to the same label with a tiny probability over the choice of the edge  $e = (u, v)$ . Technically, we express the error term in the invariance principle as a certain fourth moment of the coefficients of the halfspace, and use the smoothness to bound this error term for most edges of the Label Cover instance.

## 2.3 Bypassing the Unique Games conjecture

Unlike previous invariance principle based proofs which are only known to give Unique-Games hardness, we are able to reduce from a version of the Label Cover problem, based on *unique on average* projections, that can be shown to be NP-hard. It is of great interest to find other applications where a *weak uniqueness* property like the smoothness condition mentioned above can be used to convert a Unique-Games hardness result to an unconditional NP-hardness result. Indeed, inspired by the success of this work in avoiding the UGC assumption and using some of our methods, follow-up work has managed to bypass the Unique Games conjecture in some optimal geometric inapproximability results [?]. To the best of our knowledge, the results of [?] are the first NP-hardness proofs showing a tight inapproximability factor that is related to fundamental

parameters of Gaussian space, and among the small handful of results where optimality of a non-trivial semidefinite programming based algorithm is shown under the assumption  $P \neq NP$ . We hope that this paper has thus opened the avenue to convert at least some of the many tight Unique-Games hardness results to NP-hardness results.

### 3 Preliminaries

In this section, we define two important tools in our analysis: i) critical index, ii) invariance principle.

#### 3.1 Critical Index

The notion of critical index was first introduced by Servedio [41] and plays an important role in the analysis of halfspaces in [35, 38, 13].

**Definition 3.1.** *Given any real vector  $\mathbf{w} = (w^{(1)}, w^{(2)}, \dots, w^{(n)}) \in \mathbb{R}^n$ . Reorder the coordinates by decreasing absolute value, i.e.,  $|w^{(i_1)}| \geq |w^{(i_2)}| \geq \dots \geq |w^{(i_n)}|$  and denote  $\sigma_t^2 = \sum_{j=t}^n |w^{(i_j)}|^2$ . For  $0 \leq \tau \leq 1$ , the  $\tau$ -critical index of the vector  $\mathbf{w}$  is defined to be the smallest index  $k$  such  $|w^{(i_k)}| \leq \tau \sigma_k$ . If no such  $k$  exists ( $\forall k, |w^{(i_k)}| > \tau \sigma_k$ ), the  $\tau$ -critical index is defined to be  $+\infty$ . The vector  $\mathbf{w}$  is said to be  $\tau$ -regular if the  $\tau$ -critical index is 1.*

A simple observation from [13] is that if the critical index of a sequence is large then the sequence must contain a geometrically decreasing subsequence.

**Lemma 3.2.** *(Lemma 5.5 in [13]) Given a vector  $\mathbf{w} = (w^{(i)})_{i=1}^n$  such that  $|w^{(1)}| \geq |w^{(2)}| \geq \dots \geq |w^{(n)}|$ , if the  $\tau$ -critical index of the vector  $\mathbf{w}$  is larger than  $l$ , then for any  $1 \leq i \leq j \leq l + 1$ ,*

$$|w^{(j)}| \leq \sigma_j \leq (\sqrt{1 - \tau^2})^{j-i} \sigma_i \leq (\sqrt{1 - \tau^2})^{j-i} |w^{(i)}| / \tau.$$

*In particular, if  $j > i + (4/\tau^2) \ln(1/\tau)$  then  $|w^{(j)}| \leq |w^{(i)}|/3$ .*

For a  $\tau$ -regular weight vector, the following lemma bounds the probability that its weighted sum falls into a small interval under certain distributions on the points. The proof is in Appendix B.

**Lemma 3.3.** *Let  $\mathbf{w} \in \mathbb{R}^n$  be a  $\tau$ -regular vector  $\mathbf{w}$ , and  $\sum |w^{(i)}|^2 = 1$ .  $\mathcal{D}$  is a distribution over  $\{0, 1\}^n$ . Define a distribution  $\tilde{\mathcal{D}}$  on  $\{0, 1\}^n$  as follows: to generate  $\mathbf{y}$  from  $\tilde{\mathcal{D}}$ , first sample  $\mathbf{x}$  from  $\mathcal{D}$  and then define,*

$$y^{(i)} = \begin{cases} x^{(i)} & \text{with probability } 1 - \gamma \\ \text{random bit} & \text{with probability } \gamma. \end{cases}$$

*Then for any interval  $[a, b]$ , we have*

$$\Pr \left[ \langle \mathbf{w}, \mathbf{y} \rangle \in [a, b] \right] \leq \frac{4|b - a|}{\sqrt{\gamma}} + \frac{4\tau}{\sqrt{\gamma}} + 2e^{-\frac{\gamma^2}{2\tau^2}}.$$

Intuitively, by the Berry-Esseen Theorem,  $\langle \mathbf{w}, \mathbf{y} \rangle$  is  $\tau$  close to the Gaussian distribution if each  $y^{(i)}$  is a random bit; therefore we can bound the probability that  $\langle \mathbf{w}, \mathbf{y} \rangle$  falls into the interval  $[a, b]$ . In above lemma, each  $y^{(i)}$  has probability  $\gamma$  to be a random bit, then  $\gamma$  fraction of  $y^{(i)}$  is set to be a random bit and we can similarly bound the probability that  $\langle \mathbf{w}, \mathbf{y} \rangle$  falls into the interval  $[a, b]$ .

**Definition 3.4.** For a vector  $\mathbf{w} \in \mathbb{R}^n$ , define set of indices  $H_t(\mathbf{w}) \subseteq [n]$  as the set of indices containing the  $t$  biggest coordinates of  $\mathbf{w}$  by absolute value. Suppose its  $\tau$ -critical index is  $c_\tau$ , define set of indices  $C_\tau(\mathbf{w}) = H_{c_\tau}(\mathbf{w})$ . In other words,  $C_\tau(\mathbf{w})$  is the set of indices whose deletion makes the vector  $w$  to be  $\tau$ -regular.

**Definition 3.5.** For a vector  $\mathbf{w} \in \mathbb{R}^n$  and a subset of indices  $S \subseteq [n]$ , define the vector  $\text{Truncate}(\mathbf{w}, S) \in \mathbb{R}^n$  as:

$$(\text{Truncate}(\mathbf{w}, S))^{(i)} = \begin{cases} w^{(i)} & \text{if } i \in S \\ 0 & \text{otherwise} \end{cases}$$

As suggested by Lemma 3.2, a weight vector with a large critical index has a geometrically decreasing subsequence. The following two lemmas use this fact to bound the probability that the weighted sum of a geometrically decreasing sequence of weights falls into a small interval. First, we restate Claim 5.7 from [13] here.

**Lemma 3.6.** [Claim 5.7, [13]] Let  $\mathbf{w} = (w^{(1)}, \dots, w^{(T)})$  be such that  $|w^{(1)}| \geq |w^{(2)}| \dots \geq |w^{(T)}| \geq 0$  and  $|w^{(i+1)}| \leq \frac{|w^{(i)}|}{3}$  for  $1 \leq i \leq T-1$ . Then for any interval  $I = [\alpha - \frac{w^{(T)}}{6}, \alpha + \frac{w^{(T)}}{6}]$  of length  $\frac{|w^{(T)}|}{3}$ , there is at most one point  $\mathbf{x} \in \{0, 1\}^T$  such that  $\langle \mathbf{w}, \mathbf{x} \rangle \in I$ .

**Lemma 3.7.** Let  $\mathbf{w} = (w^{(1)}, \dots, w^{(T)})$  be such that  $|w^{(1)}| \geq |w^{(2)}| \dots \geq |w^{(T)}| \geq 0$  and  $|w^{(i+1)}| \leq \frac{|w^{(i)}|}{3}$  for  $1 \leq i \leq T-1$ . Let  $\mathcal{D}$  be a distribution over  $\{0, 1\}^T$ . Define a distribution  $\tilde{\mathcal{D}}$  on  $\{0, 1\}^T$  as follows: To generate  $\mathbf{y}$  from  $\tilde{\mathcal{D}}$ , sample  $\mathbf{x}$  from  $\mathcal{D}$  and set

$$y^{(i)} = \begin{cases} x^{(i)} & \text{with probability } 1 - \gamma \\ \text{random bit} & \text{with probability } \gamma. \end{cases}$$

Then for any  $\theta \in \mathbb{R}$  we have

$$\Pr \left[ \langle \mathbf{w}, \mathbf{y} \rangle \in \left[ \theta - \frac{w^{(T)}}{6}, \theta + \frac{w^{(T)}}{6} \right] \right] \leq \left( 1 - \frac{\gamma}{2} \right)^T.$$

*Proof.* By Lemma 3.6, we know that for the interval  $J = \left[ \theta - \frac{|w^{(T)}|}{6}, \theta + \frac{|w^{(T)}|}{6} \right]$ , there is at most one point  $\mathbf{r} \in \{0, 1\}^T$  such that  $\langle \mathbf{w}, \mathbf{r} \rangle \in J$ . If no such  $\mathbf{r}$  exists then clearly the probability is zero. On the other hand, suppose there exists such an  $\mathbf{r}$ , then  $\langle \mathbf{w}, \mathbf{y} \rangle \in J$  only if  $(y_1^{(1)}, y_1^{(2)}, \dots, y_1^{(T)}) = (r^{(1)}, \dots, r^{(T)})$  holds.

Conditioned on any fixing of the bits  $\mathbf{x}$ , every bit  $y^{(j)}$  is an independent random bit with probability  $\gamma$ . Therefore, for every fixing of  $\mathbf{x}$ , for each  $i \in [T]$ , with probability at least  $\gamma/2$ ,  $y^{(i)}$  is not equal to  $r^{(i)}$ . Therefore,  $\Pr[y^{(1)} = r^{(1)}, y^{(2)} = r^{(2)}, \dots, y^{(T)} = r^{(T)}] \leq \left( 1 - \frac{\gamma}{2} \right)^T$ .  $\square$

### 3.2 Invariance Principle

While invariance principles have been shown in various settings by [37, 11, 36], we restate a version of the principle well suited for our application. We present a self-contained proof for it in Appendix C.

**Definition 3.8.** A function  $\Psi(x) : \mathbb{R} \rightarrow \mathbb{R}$  for which fourth-order derivatives exist everywhere on  $\mathbb{R}$  is said to be  $K$ -bounded if  $|\Psi''''(t)| \leq K$  for all  $t \in \mathbb{R}$ .

**Definition 3.9.** Two ensembles of random variables  $\mathcal{P} = (p_1, \dots, p_k)$  and  $\mathcal{Q} = (q_1, \dots, q_k)$  are said to have matching moments up to degree  $d$  if for every multi-set  $S$  of elements from  $[k]$ ,  $|S| \leq d$ , we have  $\mathbf{E}[\prod_{i \in S} p_i] = \mathbf{E}[\prod_{i \in S} q_i]$ .

**Theorem 3.10.** (*Invariance Principle*) Let  $\mathcal{A} = \{\mathbf{A}^{\{1\}}, \dots, \mathbf{A}^{\{R\}}\}, \mathcal{B} = \{\mathbf{B}^{\{1\}}, \dots, \mathbf{B}^{\{R\}}\}$  be families of ensembles of random variables with  $\mathbf{A}^{\{i\}} = \{a_1^{(i)}, \dots, a_{k_i}^{(i)}\}$  and  $\mathbf{B}^{\{i\}} = \{b_1^{(i)}, \dots, b_{k_i}^{(i)}\}$ , satisfying the following properties:

- For each  $i \in [R]$ , the random variables in ensembles  $(\mathbf{A}^{\{i\}}, \mathbf{B}^{\{i\}})$  have matching moments up to degree 3. Further all the random variables in  $\mathcal{A}$  and  $\mathcal{B}$  are bounded by 1.
- The ensembles  $\mathbf{A}^{\{i\}}$  are all independent of each other, similarly the ensembles  $\mathbf{B}^{\{i\}}$  are independent of each other.

Given a set of vectors  $\mathbf{l} = \{\mathbf{l}^{\{1\}}, \dots, \mathbf{l}^{\{R\}}\} (\mathbf{l}^{\{i\}} \in \mathbb{R}^{k_i})$ , define the linear function  $\mathbf{l} : \mathbb{R}^{k_1} \times \dots \times \mathbb{R}^{k_R} \rightarrow \mathbb{R}$  as

$$\mathbf{l}(\mathbf{x}) = \sum_{i \in [R]} \langle \mathbf{l}^{\{i\}}, \mathbf{x}^{\{i\}} \rangle$$

Then for a  $K$ -bounded function  $\Psi : \mathbb{R} \rightarrow \mathbb{R}$  we have

$$\left| \mathbf{E}_{\mathcal{A}} \left[ \Psi(\mathbf{l}(\mathcal{A}) - \theta) \right] - \mathbf{E}_{\mathcal{B}} \left[ \Psi(\mathbf{l}(\mathcal{B}) - \theta) \right] \right| \leq K \sum_{i \in [R]} \|\mathbf{l}^{\{i\}}\|_1^4$$

for all  $\theta > 0$ . Further, define the spread function  $c(\alpha)$  corresponding to the ensembles  $\mathcal{A}, \mathcal{B}$  and the linear function  $\mathbf{l}$  as follows,

(**Spread Function:**) For  $1/2 > \alpha > 0$ , let

$$c(\alpha) = \max \left( \sup_{\theta} \Pr_{\mathcal{A}} \left[ \mathbf{l}(\mathcal{A}) \in [\theta - \alpha, \theta + \alpha] \right], \sup_{\theta} \Pr_{\mathcal{B}} \left[ \mathbf{l}(\mathcal{B}) \in [\theta - \alpha, \theta + \alpha] \right] \right)$$

then for all  $\theta$ ,

$$\left| \mathbf{E}_{\mathcal{A}} [\text{pos}(\mathbf{l}(\mathcal{A}) - \theta)] - \mathbf{E}_{\mathcal{B}} [\text{pos}(\mathbf{l}(\mathcal{B}) - \theta)] \right| \leq O \left( \frac{1}{\alpha^4} \right) \sum_{i \in [R]} \|\mathbf{l}^{\{i\}}\|_1^4 + 2c(\alpha).$$

Roughly speaking, the second part of the theorem states that pos function can be thought of as  $\frac{1}{\alpha^4}$ -bounded with error parameter  $c(\alpha)$ .

## 4 Construction of the Dictatorship Test

In this section we describe the construction of the dictatorship test which will be the key ingredient in the hardness reduction from  $k$ -UNIQUE LABEL COVER.

### 4.1 Distributions $\mathcal{D}_0$ and $\mathcal{D}_1$

The dictatorship test is based on following two distributions  $\mathcal{D}_0$  and  $\mathcal{D}_1$  defined on  $\{0, 1\}^k$ .

**Lemma 4.1.** *For  $k \in \mathbb{N}$ , there exists two probability distributions  $\mathcal{D}_0, \mathcal{D}_1$  on  $\{0, 1\}^k$  such that for  $x = (x_1, \dots, x_k)$ ,*

$$\Pr_{x \sim \mathcal{D}_0} \{\text{every } x_l \text{ is } 0\} \geq 1 - \frac{2}{\sqrt{k}} \text{ and } \Pr_{x \sim \mathcal{D}_1} \{\text{every } x_l \text{ is } 0\} \leq \frac{1}{\sqrt{k}},$$

while matching moments up to degree 4, i.e.,  $\forall i, j, m, n \in [k]$

$$\begin{aligned} \mathbf{E}_{\mathcal{D}_0}[x_i] &= \mathbf{E}_{\mathcal{D}_1}[x_i] & \mathbf{E}_{\mathcal{D}_0}[x_i x_j x_m x_n] &= \mathbf{E}_{\mathcal{D}_1}[x_i x_j x_m x_n] \\ \mathbf{E}_{\mathcal{D}_0}[x_i x_j] &= \mathbf{E}_{\mathcal{D}_1}[x_i x_j] & \mathbf{E}_{\mathcal{D}_0}[x_i x_j x_m] &= \mathbf{E}_{\mathcal{D}_1}[x_i x_j x_m] \end{aligned}$$

*Proof.* For  $\epsilon = \frac{1}{\sqrt{k}}$ , take  $\mathcal{D}_1$  to be the following distribution:

1. with probability  $(1 - \epsilon)$ , randomly set exactly one of the bit to be 1 and all the other to be 0;
2. with probability  $\frac{\epsilon}{4}$ , independently set every bit to be 1 with probability  $\frac{1}{k^{1/3}}$ ;
3. with probability  $\frac{\epsilon}{4}$ , independently set every bit to be 1 with probability  $\frac{2}{k^{1/3}}$ ;
4. with probability  $\frac{\epsilon}{4}$ , independently set every bit to be 1 with probability  $\frac{3}{k^{1/3}}$ ;
5. with probability  $\frac{\epsilon}{4}$ , independently set every bit to be 1 with probability  $\frac{4}{k^{1/3}}$ .

The distribution  $\mathcal{D}_0$  is defined to be the following distribution with parameter  $\epsilon_1, \epsilon_2, \epsilon_3, \epsilon_4$  to be specified later:

1. with probability  $1 - (\epsilon_1 + \epsilon_2 + \epsilon_3 + \epsilon_4)$ , set every bit to be zero;
2. with probability  $\epsilon_1$ , independently set every bit to be 1 with probability  $\frac{1}{k^{1/3}}$ ;
3. with probability  $\epsilon_2$ , independently set every bit to be 1 with probability  $\frac{2}{k^{1/3}}$ ;
4. with probability  $\epsilon_3$ , independently set every bit to be 1 with probability  $\frac{3}{k^{1/3}}$ ;
5. with probability  $\epsilon_4$ , independently set every bit to be 1 with probability  $\frac{4}{k^{1/3}}$ .

From the definition of  $\mathcal{D}_0, \mathcal{D}_1$ , we know that  $\Pr_{x \sim \mathcal{D}_0}[\text{every } x_i \text{ is } 0] \geq 1 - (\epsilon_1 + \epsilon_2 + \epsilon_3 + \epsilon_4)$  and  $\Pr_{x \sim \mathcal{D}_1}[\text{every } x_i \text{ is } 0] \leq \epsilon = \frac{1}{\sqrt{k}}$ .

It remains to determine each  $\epsilon_i$ . Notice that the moment matching conditions can be expressed as a linear system over the parameters  $\epsilon_1, \epsilon_2, \epsilon_3, \epsilon_4$  as follows:

$$\begin{aligned} \sum_{i=1}^4 \epsilon_i \left(\frac{i}{k^{1/3}}\right) &= (1 - \epsilon)/k + \sum_{i=1}^4 \frac{\epsilon}{4} \left(\frac{i}{k^{1/3}}\right) \\ \sum_{i=1}^4 \epsilon_i \left(\frac{i}{k^{1/3}}\right)^2 &= \sum_{i=1}^4 \frac{\epsilon}{4} \left(\frac{i}{k^{1/3}}\right)^2 \\ \sum_{i=1}^4 \epsilon_i \left(\frac{i}{k^{1/3}}\right)^3 &= \sum_{i=1}^4 \frac{\epsilon}{4} \left(\frac{i}{k^{1/3}}\right)^3 \\ \sum_{i=1}^4 \epsilon_i \left(\frac{i}{k^{1/3}}\right)^4 &= \sum_{i=1}^4 \frac{\epsilon}{4} \left(\frac{i}{k^{1/3}}\right)^4. \end{aligned}$$

We then show that such a linear system has a feasible solution  $\epsilon_1, \epsilon_2, \epsilon_3, \epsilon_4 > 0$  and  $\sum_{i=1}^4 \epsilon_i \leq 2/\sqrt{k}$ .

To prove this, by applying Cramer's rule,

$$\epsilon_1 = \frac{\begin{vmatrix} (1 - \epsilon)/k + \sum_{i=1}^4 \frac{\epsilon}{4} \left(\frac{i}{k^{1/3}}\right) & \frac{2}{k^{1/3}} & \frac{3}{k^{1/3}} & \frac{4}{k^{1/3}} \\ \sum_{i=1}^4 \frac{\epsilon}{4} \left(\frac{i}{k^{1/3}}\right)^2 & \frac{4}{k^{2/3}} & \frac{9}{k^{2/3}} & \frac{16}{k^{2/3}} \\ \sum_{i=1}^4 \frac{\epsilon}{4} \left(\frac{i}{k^{1/3}}\right)^3 & \frac{8}{k^{3/3}} & \frac{27}{k^{3/3}} & \frac{64}{k^{3/3}} \\ \sum_{i=1}^4 \frac{\epsilon}{4} \left(\frac{i}{k^{1/3}}\right)^4 & \frac{16}{k^{4/3}} & \frac{81}{k^{4/3}} & \frac{256}{k^{4/3}} \end{vmatrix}}{\begin{vmatrix} \frac{1}{k^{1/3}} & \frac{2}{k^{1/3}} & \frac{3}{k^{1/3}} & \frac{4}{k^{1/3}} \\ \frac{1}{k^{2/3}} & \frac{4}{k^{2/3}} & \frac{9}{k^{2/3}} & \frac{16}{k^{2/3}} \\ \frac{1}{k^{3/3}} & \frac{8}{k^{3/3}} & \frac{27}{k^{3/3}} & \frac{64}{k^{3/3}} \\ \frac{1}{k^{4/3}} & \frac{16}{k^{4/3}} & \frac{81}{k^{4/3}} & \frac{256}{k^{4/3}} \end{vmatrix}}$$

With some calculation using basic linear algebra, we get

$$\epsilon_1 = \epsilon/4 + \frac{\begin{vmatrix} (1 - \epsilon)/k & \frac{2}{k^{1/3}} & \frac{3}{k^{1/3}} & \frac{4}{k^{1/3}} \\ 0 & \frac{4}{k^{2/3}} & \frac{9}{k^{2/3}} & \frac{16}{k^{2/3}} \\ 0 & \frac{8}{k^{3/3}} & \frac{27}{k^{3/3}} & \frac{64}{k^{3/3}} \\ 0 & \frac{16}{k^{4/3}} & \frac{81}{k^{4/3}} & \frac{256}{k^{4/3}} \end{vmatrix}}{\begin{vmatrix} \frac{1}{k^{1/3}} & \frac{2}{k^{1/3}} & \frac{3}{k^{1/3}} & \frac{4}{k^{1/3}} \\ \frac{1}{k^{2/3}} & \frac{4}{k^{2/3}} & \frac{9}{k^{2/3}} & \frac{16}{k^{2/3}} \\ \frac{1}{k^{3/3}} & \frac{8}{k^{3/3}} & \frac{27}{k^{3/3}} & \frac{64}{k^{3/3}} \\ \frac{1}{k^{4/3}} & \frac{16}{k^{4/3}} & \frac{81}{k^{4/3}} & \frac{256}{k^{4/3}} \end{vmatrix}} = \frac{1}{4\sqrt{k}} + O\left(\frac{1}{k^{2/3}}\right).$$

For large enough  $k$ , we have  $0 \leq \epsilon_1 \leq \frac{1}{2\sqrt{k}}$ . By similar calculation, we can bound  $\epsilon_2, \epsilon_3, \epsilon_4$  by  $\frac{1}{2\sqrt{k}}$ . Overall, we have  $\epsilon_1 + \epsilon_2 + \epsilon_3 + \epsilon_4 \leq 2/\sqrt{k}$

□

We define a “noisy” version of  $\mathcal{D}_b$  ( $b \in \{0, 1\}$ ) below.

**Definition 4.2.** For  $b \in \{0, 1\}$ , define the distribution  $\tilde{\mathcal{D}}_b$  on  $\{0, 1\}^k$  as follows:

- First generate  $x \in \{0, 1\}^k$  according to  $\mathcal{D}_b$ .
- For each  $i \in [k]$ ,

$$y_i = \begin{cases} x_i & \text{with probability } 1 - \frac{1}{k^2} \\ \text{uniform random bit } u_i & \text{with probability } \frac{1}{k^2} \end{cases}$$

**Observation 4.3.**  $\tilde{\mathcal{D}}_0$  and  $\tilde{\mathcal{D}}_1$  also have matching moments up to degree 4.

*Proof.* Since the noise is defined to be an independent uniform random bit, when calculating moments of  $y$ , such as  $\mathbf{E}_{\tilde{\mathcal{D}}_b}[y_{i_1}y_{i_2} \cdots y_{i_d}]$ , we can substitute  $y_i$  by  $(1 - \frac{1}{k^2})x_i + \frac{1}{k^2}$ . Therefore, a degree  $d$  moment of  $y$  can be expressed as a weighted sum of moments of  $x$  of degree up to  $d$ . Since  $\mathcal{D}_0$  and  $\mathcal{D}_1$  have matching moments up to degree 4, it follows that  $\tilde{\mathcal{D}}_0$  and  $\tilde{\mathcal{D}}_1$  also have the same property. □

The following simple lemma asserts that conditioning the two distributions  $\tilde{\mathcal{D}}_0$  and  $\tilde{\mathcal{D}}_1$  on the same coordinate  $x_j$  being fixed to value  $b$  results in conditional distributions that still have matching moments up to degree 3.

**Lemma 4.4.** Given two distributions  $\mathcal{P}_0, \mathcal{P}_1$  on  $\{0, 1\}^k$  with matching moments up to degree  $d$ , for any multi-set  $S$  of elements from  $[k]$ ,  $|S| \leq d - 1$ ,  $j \in [k]$  and  $c \in \{0, 1\}$ .

$$\mathbf{E}_{\mathcal{P}_0}[\prod_{i \in S} x_i \mid x_j = c] = \mathbf{E}_{\mathcal{P}_1}[\prod_{i \in S} x_i \mid x_j = c].$$

*Proof.* For the case  $c = 1$  and any  $b \in \{0, 1\}$ ,

$$\mathbf{E}_{\mathcal{P}_b}[x_j \prod_{i \in S} x_i] = \mathbf{E}_{\mathcal{P}_b}[\prod_{i \in S} x_i \mid x_j = 1] \Pr_{\mathcal{P}_0}[x_j = 1] = \mathbf{E}_{\mathcal{P}_b}[\prod_{i \in S} x_i \mid x_j = 1] \mathbf{E}_{\mathcal{P}_0}[x_j].$$

Therefore,

$$\mathbf{E}_{\mathcal{P}_0}[\prod_{i \in S} x_i \mid x_j = 1] = \frac{\mathbf{E}_{\mathcal{P}_0}[x_j \prod_{i \in S} x_i]}{\mathbf{E}_{\mathcal{P}_0}[x_j]} = \frac{\mathbf{E}_{\mathcal{P}_1}[x_j \prod_{i \in S} x_i]}{\mathbf{E}_{\mathcal{P}_1}[x_j]} = \mathbf{E}_{\mathcal{P}_1}[\prod_{i \in S} x_i \mid x_j = 1].$$

For the case  $c = 0$ , replace  $x_j$  with  $x'_j = 1 - x_j$ . It is easy to see that  $\mathcal{P}_0$  and  $\mathcal{P}_1$  still have matching moments and conditioning on  $x_j = 0$  is the same as conditioning on  $x'_j = 1$ . Hence we can reduce to the case  $c = 1$ . □

## 4.2 The Dictatorship Test

Let  $R$  be a positive integer. Based on the distributions  $\mathcal{D}_0$  and  $\mathcal{D}_1$ , we define the dictatorship test as follows:

1. Generate a random bit  $b \in \{0, 1\}$ .
2. Generate  $\mathbf{x} \in \{0, 1\}^{kR}$  (which is also written as  $\{x_i^{(j)}\}_{i \in [k], j \in [R]}$ ) from  $\mathcal{D}_b^R$ .
3. For each  $i \in [k], j \in [R]$ ,

$$y_i^{(j)} = \begin{cases} x_i^{(j)} & \text{with probability } 1 - \frac{1}{k^2}; \\ \text{random bit} & \text{with probability } \frac{1}{k^2}. \end{cases}$$

4. Output the labelled example  $(\mathbf{y}, b)$ . Equivalently, if  $h$  denotes the halfspace, ACCEPT if  $h(\mathbf{y}) = b$ .

We can also view  $y$  as being generated as follows: i) With probability  $\frac{1}{2}$ , generate a negative sample from distribution  $\tilde{\mathcal{D}}_0^R$ ; ii) With probability  $\frac{1}{2}$ , generate a positive sample from distribution  $\tilde{\mathcal{D}}_1^R$ .

The dictatorship test has the following completeness and soundness properties.

**Theorem 4.5.** (completeness) For any  $j \in [R]$ ,  $h(\mathbf{y}) = \bigvee_{i=1}^k y_i^{(j)}$  passes with probability  $\geq 1 - \frac{3}{\sqrt{k}}$ .

**Theorem 4.6.** (soundness) Fix  $\tau = \frac{1}{k^7}$  and  $t = \frac{1}{\tau^2}(3 \ln(1/\tau) + \ln R) + \lceil 4k^2 \ln k \rceil \lceil \frac{4}{\tau^2} \ln(1/\tau) \rceil$ . Let  $h(\mathbf{x}) = \text{pos}(\langle \mathbf{w}, \mathbf{y} \rangle - \theta)$  be a halfspace such that  $H_t(\mathbf{w}_i) \cap H_t(\mathbf{w}_j) = \emptyset$  for all  $i, j \in [k]$ . Then the halfspace  $h(\mathbf{y})$  passes the dictatorship test with probability at most  $\frac{1}{2} + O(\frac{1}{k})$ .

*Proof.* (Theorem 4.5) If  $x$  is generated from  $\mathcal{D}_0^R$ , we know that with probability at least  $1 - \frac{2}{\sqrt{k}}$ , all the bits in  $\{x_1^{(j)}, x_2^{(j)}, \dots, x_k^{(j)}\}$  are set to 0. By union bound, with probability at least  $1 - \frac{2}{\sqrt{k}} - \frac{1}{k}$ ,  $\{y_1^{(j)}, y_2^{(j)}, \dots, y_k^{(j)}\}$  are all set to 0, in which case the test passes as  $\bigvee_{i=1}^k y_i^{(j)} = 0$ . If  $x$  is generated from  $\mathcal{D}_1^R$ , we know that with probability at least  $1 - \frac{1}{\sqrt{k}}$ , one of the bits in  $\{x_1^{(j)}, x_2^{(j)}, \dots, x_k^{(j)}\}$  is set to 1 and by union bound one of  $\{y_1^{(j)}, y_2^{(j)}, \dots, y_k^{(j)}\}$  is set to 1 with probability at least  $1 - \frac{1}{\sqrt{k}} - \frac{1}{k}$ , in which case the test passes since  $\bigvee_{i=1}^k y_i^{(j)} = 1$ . Overall, the test passes with probability at least  $1 - \frac{3}{\sqrt{k}}$ .  $\square$

## 4.3 Proof of Soundness (Theorem 4.6)

We will prove the contrapositive statement of Theorem 4.6: if some  $h(\mathbf{y})$  passes the above dictatorship test with high probability, then we can decode for each  $\mathbf{w}_i$  ( $i \in [k]$ ), a small list of coordinates and at least two of the lists will intersect.

The proof is based on two key lemmas (Lemmas 4.7, 4.8). The first lemma states that if a halfspace passes the test with good probability, then two of its critical index sets  $C_\tau(\mathbf{w}_i), C_\tau(\mathbf{w}_j)$

must intersect. This would immediately imply Theorem 4.6 if  $c_\tau$  is less than  $t$ . The second lemma states that every halfspace can be approximated by another halfspace with critical index less than  $t$ ; so we can assume that  $c_\tau$  is small without loss of generality.

Let  $h(\mathbf{y})$  be a halfspace function on  $\{0, 1\}^{kR}$  given by  $h(\mathbf{y}) = \text{pos}(\langle \mathbf{w}, \mathbf{y} \rangle - \theta)$ . Equivalently,  $h(\mathbf{y})$  can be written as

$$h(\mathbf{y}) = \text{pos}\left(\sum_{j \in [R]} \langle \mathbf{w}^{(j)}, \mathbf{y}^{(j)} \rangle - \theta\right) = \text{pos}\left(\sum_{i \in [k]} \langle \mathbf{w}_i, \mathbf{y}_i \rangle - \theta\right),$$

where  $\mathbf{w}^{(j)} \in \mathbb{R}^k$  and  $\mathbf{w}_i \in \mathbb{R}^R$ .

**Lemma 4.7.** (*Common Influential Coordinates*) For  $\tau = \frac{1}{k^\tau}$ , let  $h(\mathbf{y})$  be a halfspace such that for all  $i \neq j \in [k]$ , we have  $C_\tau(\mathbf{w}_i) \cap C_\tau(\mathbf{w}_j) = \emptyset$ . Then

$$\left| \mathbf{E}_{\tilde{\mathcal{D}}_0^R} [h(\mathbf{y})] - \mathbf{E}_{\tilde{\mathcal{D}}_1^R} [h(\mathbf{y})] \right| \leq O\left(\frac{1}{k}\right).$$

*Proof.* Fix the following notation,

$$\begin{aligned} \mathbf{l}_i &= \text{Truncate}(\mathbf{w}_i, C_\tau(\mathbf{w}_i)) & \mathbf{s}_i &= \mathbf{w}_i - \mathbf{w}_i^C \\ \mathbf{y}_i^C &= \text{Truncate}(\mathbf{y}_i, C_\tau(\mathbf{w}_i)) & \mathbf{y}^C &= \mathbf{y}_1^C, \mathbf{y}_2^C, \dots, \mathbf{y}_k^C \\ \mathbf{s} &= \mathbf{s}_1, \mathbf{s}_2, \dots, \mathbf{s}_k & \mathbf{l} &= \mathbf{l}_1, \mathbf{l}_2, \dots, \mathbf{l}_k. \end{aligned}$$

We can rewrite the halfspace  $h(\mathbf{y})$  as  $h(\mathbf{y}) = \text{pos}(\langle \mathbf{l}, \mathbf{y}^C \rangle + \langle \mathbf{s}, \mathbf{y} \rangle - \theta)$ . Let us first normalize the halfspace  $h(\mathbf{y})$  so that  $\sum_{i \in [k]} \|\mathbf{l}_i\|^2 = 1$ . We now condition on a possible fixing of the vector  $\mathbf{y}^C$ . Under this conditioning and for  $\mathbf{y}$  chosen randomly from the distribution  $\tilde{\mathcal{D}}_0^R$ , define the family of ensembles  $\mathcal{A} = \mathbf{A}^{\{1\}}, \dots, \mathbf{A}^{\{R\}}$  as follows:

$$\mathbf{A}^{\{j\}} = \{y_i^{(j)} \mid i \in [k] \text{ for which } j \notin C_\tau(\mathbf{w}_i)\}$$

Similarly define the ensemble  $\mathcal{B} = \mathbf{B}^{\{1\}}, \dots, \mathbf{B}^{\{R\}}$  using  $\mathbf{y}$  chosen randomly from the distribution  $\tilde{\mathcal{D}}_1^R$ . Further let us denote  $\mathbf{l}^{\{j\}} = (l_1^{(j)}, \dots, l_k^{(j)})$ . Now we apply the invariance principle (Theorem 3.10) to the ensembles  $\mathcal{A}, \mathcal{B}$  and the linear function  $\mathbf{l}$ . For each  $j \in [R]$ , there is at most one coordinate  $i \in [k]$  such that  $j \in C_\tau(\mathbf{w}_i)$ . Thus, conditioning on  $\mathbf{y}^C$  amounts to fixing of at most one variable  $y_i^{(j)}$  in each column  $\{y_i^{(j)}\}_{i \in [k]}$ . By Lemma 4.4, since  $\tilde{\mathcal{D}}_0$  and  $\tilde{\mathcal{D}}_1$  have matching moments up to degree 4, we get that  $\mathbf{A}^{\{j\}}$  and  $\mathbf{B}^{\{j\}}$  have matching moments up to degree 3. Also notice that  $\max_{j \in [R], i \in [k]} |l_i^{(j)}| \leq \tau \|\mathbf{l}_i\|_2 \leq \tau \|\mathbf{l}\|_2$  (as  $\mathbf{l}_i$  is a  $\tau$ -regular) and each  $y_i^{(j)}$  is set to be a random unbiased bit with probability  $\frac{1}{k^2}$ ; by Lemma 3.3, the linear function  $\mathbf{l}$  and the ensembles  $\mathcal{A}, \mathcal{B}$  satisfy the following spread property for every  $\theta' \in \mathbb{R}$ :

$$\begin{aligned} \Pr_{\mathcal{A}} \left[ \mathbf{l}(\mathcal{A}) \in [\theta' - \alpha, \theta' + \alpha] \right] &\leq c(\alpha) \\ \Pr_{\mathcal{B}} \left[ \mathbf{l}(\mathcal{B}) \in [\theta' - \alpha, \theta' + \alpha] \right] &\leq c(\alpha), \end{aligned}$$

where  $c(\alpha) \leq 8\alpha k + 4\tau k + 2e^{-\frac{1}{2\tau^2 k^4}}$  (by setting  $\gamma = \frac{1}{k^2}$  and  $|b - a| = 2\alpha$  in Lemma 3.3). Using the invariance principle (Theorem 3.10) this implies:

$$\begin{aligned} & \left| \mathbf{E}_{\mathcal{A}} \left[ \text{pos} \left( \langle \mathbf{s}, \mathbf{y}^C \rangle + \sum_{j \in [R]} \langle \mathbf{l}^{\{j\}}, \mathbf{A}^{\{j\}} \rangle - \theta \right) \middle| \mathbf{y}^C \right] - \right. \\ & \quad \left. \mathbf{E}_{\mathcal{B}} \left[ \text{pos} \left( \langle \mathbf{s}, \mathbf{y}^C \rangle + \sum_{j \in [R]} \langle \mathbf{l}^{\{j\}}, \mathbf{B}^{\{j\}} \rangle - \theta \right) \middle| \mathbf{y}^C \right] \right| \\ & \leq O\left(\frac{1}{\alpha^4}\right) \sum_{i \in [R]} \|\mathbf{l}^{\{i\}}\|_1^4 + 2c(\alpha) \quad (1) \end{aligned}$$

By definition of the critical index, we have  $\max_{j \in [R]} l_i^j \leq \tau \|\mathbf{l}_i\|_2$ . Using this, we can bound  $\sum_{i \in [R]} \|\mathbf{l}^{\{i\}}\|_1^4$  as follows:

$$\begin{aligned} \sum_{j \in [R]} \|\mathbf{l}^{\{j\}}\|_1^4 & \leq k^4 \sum_{i \in [k]} \sum_{j \in [R]} \|l_i^{(j)}\|^4 \leq k^4 \sum_{i \in [k]} \left( \max_{j \in [R]} |l_i^{(j)}|^2 \right) \|\mathbf{l}_i\|_2^2 \\ & \leq k^4 \tau^2 \sum_{i \in [k]} \|\mathbf{l}_i\|_2^2 \leq k^4 \tau^2 \|\mathbf{l}\|_2^2 \leq \frac{1}{k^{10}}. \end{aligned}$$

In the final inequality in above calculation, we used the fact that  $\tau = \frac{1}{k^7}$  and  $\|\mathbf{l}\|_2 = 1$ . Let us choose  $\alpha = \frac{1}{k^2}$  and (1) is therefore bounded by  $O(1/k)$  for all settings of  $\mathbf{y}^C$ . Averaging over all settings of  $\mathbf{y}^C$  we get that

$$\left| \mathbf{E}_{\tilde{\mathcal{D}}_0^R} [h(\mathbf{y})] - \mathbf{E}_{\tilde{\mathcal{D}}_1^R} [h(\mathbf{y})] \right| \leq O\left(\frac{1}{k}\right).$$

□

The above lemma asserts that unless some two vectors  $\mathbf{w}_i, \mathbf{w}_j$  have a *common influential coordinate*, the halfspace  $h(\mathbf{y})$  cannot distinguish between  $\tilde{\mathcal{D}}_0^R$  and  $\tilde{\mathcal{D}}_1^R$ . Unlike with the traditional notion of influence, it is unclear whether the number of coordinates in  $C_\tau(\mathbf{w}_i)$  is small. The following lemma yields a way to get around this.

**Lemma 4.8.** (*Bounding the number of influential coordinates*) *Let  $t$  be set as in Theorem 4.6. Given a halfspace  $h(\mathbf{y})$  and  $r \in [k]$  such that  $|C_\tau(\mathbf{w}_r)| > t$ , define  $\tilde{h}(\mathbf{y}) = \text{pos}(\sum_{i \in [k]} \langle \tilde{\mathbf{w}}_i, \mathbf{y}_i \rangle - \theta)$  as follows:  $\tilde{\mathbf{w}}_r = \text{Truncate}(\mathbf{w}_r, H_t(\mathbf{w}_r))$  and  $\tilde{\mathbf{w}}_i = \mathbf{w}_i$  for all  $i \neq r$ . Then,*

$$\left| \mathbf{E}_{\tilde{\mathcal{D}}_0^R} [\tilde{h}(\mathbf{y})] - \mathbf{E}_{\tilde{\mathcal{D}}_0^R} [h(\mathbf{y})] \right| \leq \frac{1}{k^2} \text{ and } \left| \mathbf{E}_{\tilde{\mathcal{D}}_1^R} [\tilde{h}(\mathbf{y})] - \mathbf{E}_{\tilde{\mathcal{D}}_1^R} [h(\mathbf{y})] \right| \leq \frac{1}{k^2}.$$

*Proof.* Without loss of generality, we assume  $r = 1$  and  $|w_1^{(1)}| \geq |w_1^{(2)}| \geq \dots \geq |w_1^{(R)}|$ . In particular, this implies  $H_t(\mathbf{w}_1) = \{1, \dots, t\}$ . Set  $T = \lceil 4k^2 \ln k \rceil$ . Define the subset  $G$  of  $H_t(\mathbf{w}_1)$  as

$$G = \{g_i \mid g_i = 1 + i \lceil (4/\tau^2) \ln(1/\tau) \rceil, 0 \leq i \leq T\}.$$

Therefore, by Lemma 3.2,  $|w_1^{(g_i)}|$  is a geometrically decreasing sequence such that  $|w_1^{(g_{i+1})}| \leq |w_1^{(g_i)}|/3$ . Let  $H = H_t(\mathbf{w}_1) \setminus G$ . Fix the following notation:

$$\mathbf{w}_1^G = \text{Truncate}(\mathbf{w}_1, G), \quad \mathbf{w}_1^H = \text{Truncate}(\mathbf{w}_1, H), \quad \mathbf{w}_1^{>t} = \text{Truncate}(\mathbf{w}_1, \{t+1, \dots, n\}).$$

Similarly, define the vectors  $\mathbf{y}_1^G, \mathbf{y}_1^H, \mathbf{y}_1^{>t}$ . We now rewrite the halfspace functions  $h(\mathbf{y})$  and  $\tilde{h}(\mathbf{y})$  as:

$$h(\mathbf{y}) = \text{pos} \left( \sum_{i=2}^k \langle \mathbf{w}_i, \mathbf{y}_i \rangle + \langle \mathbf{w}_1^G, \mathbf{y}_1^G \rangle + \langle \mathbf{w}_1^H, \mathbf{y}_1^H \rangle + \langle \mathbf{w}_1^{>t}, \mathbf{y}_1^{>t} \rangle - \theta \right)$$

$$\tilde{h}(\mathbf{y}) = \text{pos} \left( \sum_{i=2}^k \langle \mathbf{w}_i, \mathbf{y}_i \rangle + \langle \mathbf{w}_1^G, \mathbf{y}_1^G \rangle + \langle \mathbf{w}_1^H, \mathbf{y}_1^H \rangle - \theta \right).$$

Notice that for any  $\mathbf{y}$ ,  $h(\mathbf{y}) \neq \tilde{h}(\mathbf{y})$  implies

$$\left| \sum_{i=2}^k \langle \mathbf{w}_i, \mathbf{y}_i \rangle + \langle \mathbf{w}_1^G, \mathbf{y}_1^G \rangle + \langle \mathbf{w}_1^H, \mathbf{y}_1^H \rangle - \theta \right| \leq |\langle \mathbf{w}_1^{>t}, \mathbf{y}_1^{>t} \rangle|. \quad (2)$$

By Lemma 3.2, we know that

$$|w_1^{(gT)}|^2 \geq \frac{\tau^2}{(1-\tau^2)^{t-gT}} \|\mathbf{w}_1^{>t}\|_2^2 \geq \frac{\tau^2}{(1-\tau^2)^{\frac{1}{\tau^2}(3\ln(1/\tau)+\ln R)}} \|\mathbf{w}_1^{>t}\|_2^2 \geq \frac{R}{\tau} \|\mathbf{w}_1^{>t}\|_2^2.$$

Using the fact that  $R\|\mathbf{w}_1^{>t}\|_2^2 \geq \|\mathbf{w}_1^{>t}\|_1^2$ , we can get that  $\|\mathbf{w}_1^{>t}\|_1 \leq \sqrt{\tau}|w_1^{(gT)}| \leq \frac{1}{6}|w_1^{(gT)}|$ . Combining the above inequality with (2) we see that,

$$\begin{aligned} \Pr_{\tilde{\mathcal{D}}_0^R} \left[ h(\mathbf{y}) \neq \tilde{h}(\mathbf{y}) \right] &\leq \Pr_{\tilde{\mathcal{D}}_0^R} \left[ \left| \sum_{i=2}^k \langle \mathbf{w}_i, \mathbf{y}_i \rangle + \langle \mathbf{w}_1^G, \mathbf{y}_1^G \rangle + \langle \mathbf{w}_1^H, \mathbf{y}_1^H \rangle - \theta \right| \leq |\langle \mathbf{w}_1^{>t}, \mathbf{y}_1^{>t} \rangle| \right] \\ &\leq \Pr_{\tilde{\mathcal{D}}_0^R} \left[ \left| \sum_{i=2}^k \langle \mathbf{w}_i, \mathbf{y}_i \rangle + \langle \mathbf{w}_1^G, \mathbf{y}_1^G \rangle + \langle \mathbf{w}_1^H, \mathbf{y}_1^H \rangle - \theta \right| \leq \frac{|w_1^{(gT)}|}{6} \right] \\ &= \Pr_{\tilde{\mathcal{D}}_0^R} \left[ \langle \mathbf{w}_1^G, \mathbf{y}_1^G \rangle \in \left[ \theta' - \frac{1}{6}|w_1^{(gT)}|, \theta' + \frac{1}{6}|w_1^{(gT)}| \right] \right] \end{aligned}$$

where  $\theta' = -\sum_{i=2}^k \langle \mathbf{w}_i, \mathbf{y}_i \rangle - \langle \mathbf{w}_1^H, \mathbf{y}_1^H \rangle + \theta$ . For any fixing of the value of  $\theta' \in \mathbb{R}$ , it induces a certain distribution on  $\mathbf{y}_1^G$ . However, the  $\frac{1}{k^2}$  noise introduced in  $\mathbf{y}_1^G$  is completely independent. This corresponds to the setting of Lemma 3.7, and hence we can bound the above probability by  $(1 - \frac{1}{2k^2})^T \leq \frac{1}{k^2}$ . The result follows from averaging over all values of  $\theta'$ .  $\square$

With the two lemmas above, we now prove the soundness property.

*Proof.* (Theorem 4.6) The probability of success of  $h(\mathbf{y})$  is given by  $\frac{1}{2} + \frac{1}{2}(\mathbf{E}_{\tilde{\mathcal{D}}_1^R}[h(\mathbf{y})] - \mathbf{E}_{\tilde{\mathcal{D}}_0^R}[h(\mathbf{y})])$ . Therefore, it suffices to show that  $\left| \mathbf{E}_{\tilde{\mathcal{D}}_0^R}[h(\mathbf{y})] - \mathbf{E}_{\tilde{\mathcal{D}}_1^R}[h(\mathbf{y})] \right| = O(\frac{1}{k})$ .

Define  $I = \{r \mid C_\tau(\mathbf{w}_r) \geq t\}$ . We discuss the following two cases.

1.  $I = \emptyset$ ; i.e.,  $\forall i \in [k], C_\tau(\mathbf{w}_i) \leq t$ . Then for all  $i, j$ ,  $H_t(\mathbf{w}_i) \cap H_t(\mathbf{w}_j) = \emptyset$  implies  $C_\tau(\mathbf{w}_i) \cap C_\tau(\mathbf{w}_j) = \emptyset$ . By Lemma 4.7, we thus have  $\left| \mathbf{E}_{\tilde{\mathcal{D}}_0^R}[h(\mathbf{y})] - \mathbf{E}_{\tilde{\mathcal{D}}_1^R}[h(\mathbf{y})] \right| = O(\frac{1}{k})$ .

1. Sample an edge  $e = (v_1, \dots, v_k) \in E$ .
2. Generate a random bit  $b \in \{0, 1\}$ .
3. Sample  $\mathbf{x} \in \{0, 1\}^{kR}$  from  $\tilde{\mathcal{D}}_b^R$ .
4. Define  $\mathbf{y} \in \{0, 1\}^{|V| \times R}$  as follows:
  - (a) For each  $v \notin \{v_1, \dots, v_k\}$ ,  $\mathbf{y}_v = \mathbf{0}$ .
  - (b) For each  $i \in [k]$  and  $j \in [R]$ ,  $y_{v_i}^{(j)} = x_i^{(\pi^{v_i, e}(j))}$ .
5. Output the example  $(\mathbf{y}, b)$ .

Figure 1: Reduction from  $k$ -UNIQUE LABEL COVER

2.  $I \neq \emptyset$ . Then for all  $r \in I$ , we set  $\tilde{\mathbf{w}}_r = \text{Truncate}(\mathbf{w}_r, H_t(\mathbf{w}_r))$  and replace  $\mathbf{w}_r$  with  $\tilde{\mathbf{w}}_r$  in  $h$  to get a new halfspace  $h'$ . Since such replacements occur at most  $k$  times and by Lemma 4.8 every replacement changes the output of the halfspace on at most  $\frac{1}{k^2}$  fraction of examples, we can bound the overall change by  $k \times \frac{1}{k^2} = \frac{1}{k}$ . That is

$$\left| \mathbf{E}_{\tilde{\mathcal{D}}_0^R} [h'(\mathbf{y})] - \mathbf{E}_{\tilde{\mathcal{D}}_0^R} [h(\mathbf{y})] \right| \leq \frac{1}{k}, \quad \left| \mathbf{E}_{\tilde{\mathcal{D}}_1^R} [h'(\mathbf{y})] - \mathbf{E}_{\tilde{\mathcal{D}}_1^R} [h(\mathbf{y})] \right| \leq \frac{1}{k}. \quad (3)$$

Also notice that for  $h'$  and all  $r \in [k]$ , the critical index of  $\tilde{\mathbf{w}}_r$  (i.e.,  $|C_\tau(\tilde{\mathbf{w}}_r)|$ ) is less than  $t$ . This reduces the problem to Case 1, and we conclude  $\left| \mathbf{E}_{\tilde{\mathcal{D}}_0^R} [h'(\mathbf{y})] - \mathbf{E}_{\tilde{\mathcal{D}}_1^R} [h'(\mathbf{y})] \right| = O(1/k)$ . Along with (3) this finishes the proof of Theorem 4.6. □

#### 4.4 Reduction from $k$ -Unique Label Cover

With the dictatorship test defined, we now describe briefly a reduction from  $k$ -UNIQUE LABEL COVER problem to agnostic learning of monomials, thus showing Theorem 1.1 under the Unique Games Conjecture (Conjecture 2.2). Although our final hardness result only assumes  $P \neq NP$ , we describe the reduction to  $k$ -UNIQUE LABEL COVER for the purpose of illustrating the main idea of our proof.

Let  $\mathcal{L}(G(V, E), R, R, \{\pi^{v, e} | v \in V, e \in E\})$  be an instance of  $k$ -UNIQUE LABEL COVER. The reduction is defined in Figure 4.4. It will produce a distribution over labeled examples:  $(\mathbf{y}, b)$  where  $\mathbf{y} \in \{0, 1\}^{|V| \times R}$  and label  $b \in \{0, 1\}$ . We will index the coordinates of  $\mathbf{y} \in \{0, 1\}^{|V| \times R}$  by  $y_w^{(i)}$  (for  $w \in V, i \in [R]$ ) and denote  $\mathbf{y}_w$  (for  $w \in V$ ) to be the vector  $(y_w^{(1)}, y_w^{(2)}, \dots, y_w^{(R)})$ .

**Proof of Theorem 1.1 assuming Unique Games Conjecture** Fix  $k = \frac{10}{\epsilon^2}$ ,  $\eta = \frac{\epsilon^3}{100}$  and a positive integer  $R > \lceil (2k)^{\frac{1}{\eta^2}} \rceil$  for which Conjecture 2.2 holds.

**Completeness:** Suppose that  $\Lambda : V \rightarrow [R]$  is a labeling that *strongly* satisfies  $1 - k\eta$  fraction of the edges. Consider disjunction  $h(\mathbf{y}) = \bigvee_{v \in V} y_v^{(\Lambda(v))}$ . For at least  $1 - k\eta$  fraction of edges  $e = (v_1, v_2, \dots, v_k) \in E$ ,  $\pi^{v_1, e}(\Lambda(v_1)) = \dots = \pi^{v_k, e}(\Lambda(v_k)) = r$ . Let us fix such a choice of edge  $e$  in step 1. As all coordinates of  $\mathbf{y}$  outside of  $\{\mathbf{y}_{v_1}, \dots, \mathbf{y}_{v_k}\}$  are set to 0 in step 4(a), the disjunction reduces to  $\bigvee_{i \in [k]} y_{v_i}^{(\Lambda(v_i))} = \bigvee_{i \in [k]} x_i^{(r)}$ . By Theorem 4.5, such a disjunction agrees with every  $(\mathbf{y}, b)$  with probability at least  $1 - \frac{3}{\sqrt{k}}$ . Therefore  $h(\mathbf{y})$  agrees with a random example with probability at least  $(1 - \frac{3}{\sqrt{k}})(1 - k\eta) \geq 1 - \frac{3}{\sqrt{k}} - k\eta \geq 1 - \epsilon$ .

**Soundness:** Suppose there exists a halfspace  $h(\mathbf{y}) = \sum_{v \in V} \langle \mathbf{w}_v, \mathbf{y}_v \rangle$  that agrees with more than  $\frac{1}{2} + \epsilon \geq \frac{1}{2} + \frac{1}{\sqrt{k}}$  fraction of the examples. Set  $t = k^{14}(3 \ln(k^7) + \ln R) + \lceil 4k^{14} \ln k^7 \rceil \cdot \lceil 4k^2 \ln k \rceil = O(k^{16} \ln R)$  (same as in Theorem 4.6). Define the labeling  $\Lambda$  using the following strategy : for each vertex  $v \in V$  randomly pick a label from  $H_t(\mathbf{w}_v)$ .

By an averaging argument, for at least  $\frac{\epsilon}{2}$  fraction of the edges  $e \in E$  generated in step 1 of the reduction,  $h(\mathbf{y})$  agrees with the examples corresponding to  $e$  with probability at least  $\frac{1}{2} + \frac{\epsilon}{2}$ . We will refer to such edges as *good*. By Theorem 4.6 for each *good* edge  $e \in E$ , there exists  $i, j \in [k]$ , such that  $\pi^{v_i, e}(H_t(\mathbf{w}_{v_i})) \cap \pi^{v_j, e}(H_t(\mathbf{w}_{v_j})) \neq \emptyset$ . Therefore the edge  $e \in E$  is *weakly* satisfied by the labeling  $\Lambda$  with probability at least  $\frac{1}{t^2}$ . Hence, in expectation the labeling  $\Lambda$  *weakly* satisfies at least  $\frac{\epsilon}{2} \cdot \frac{1}{t^2} = \Omega(\frac{1}{k^{33} \ln^2 R}) \geq \frac{2k^2}{R^{\eta/4}}$  fraction of the edges (by the choice of  $R$  and  $t$ ).

## 5 Reduction from Label Cover

In this section, we describe a reduction from a  $k$ -LABEL COVER with an additional *smoothness* property to the problem of agnostic learning of disjunctions by halfspaces. This will give us Theorem 1.1 without assuming the Unique Games Conjecture.

### 5.1 Smooth $k$ -Label Cover

Our reduction use the following hardness result for  $k$ -LABEL COVER (Definition 2.1) with the additional smoothness property.

**Theorem 5.1.** *There exists a constant  $\gamma > 0$  such that for any integer parameter  $J, u \geq 1$ , it is NP-hard to distinguish between the following two types of  $k$ -LABEL COVER  $\mathcal{L}(G(V, E), M, N, \{\pi^{v, e} | e \in E, v \in e\})$  instances with  $M = 7^{(J+1)u}$  and  $N = 2^u 7^{Ju}$ .*

1. (Strongly satisfiable instances) *There is some labeling that strongly satisfies every hyperedge.*
2. (Instances that are not  $2k^2 2^{-\gamma u}$ -weakly satisfiable) *There is no labeling that weakly satisfies at least  $2k^2 2^{-\gamma u}$  fraction of the hyperedges.*

*In addition, the  $k$ -LABEL COVER instances have the following properties:*

- (Smoothness) *for a fixed vertex  $v$  and a randomly picked hyperedge containing  $v$ ,*

$$\forall i, j \in [M], \Pr[\pi^{v, e}(i) = \pi^{v, e}(j)] \leq 1/J.$$

- Pick a hyperedge  $e = (v_1, v_2, \dots, v_k) \in E$  with corresponding projections  $\pi^{v_1, e}, \dots, \pi^{v_k, e} : [M] \rightarrow [N]$ .
- Generate a random bit  $b \in \{0, 1\}$ .
- Sample  $\mathbf{x} \in \{0, 1\}^{kN}$  from  $\mathcal{D}_b^N$ .
- Generate  $\mathbf{y} \in \{0, 1\}^{|V| \times M}$  as follows:
  1. For each  $v \notin e$ ,  $\mathbf{y}_v = \mathbf{0}$ .
  2. For each  $i \in [k]$ , set  $\mathbf{y}_{v_i} \in \{0, 1\}^M$  as follows:
 
$$\mathbf{y}_{v_i}^{(j)} = \begin{cases} x_i^{(\pi^{v_i, e}(j))} & \text{with probability } 1 - \frac{1}{k^2} \\ \text{random bit} & \text{with probability } \frac{1}{k^2} \end{cases}$$
- Output the example  $(\mathbf{y}, b)$  or equivalently ACCEPT if  $h(\mathbf{y}) = b$ .

Figure 2: Reduction from  $k$ -LABEL COVER

- For any mapping  $\pi^{v, e}$  and any number  $i \in [N]$ , we have  $|(\pi^{v, e})^{-1}(i)| \leq d = 4^u$ ; i.e., there are at most  $d = 4^u$  elements in  $[M]$  that are mapped to the same number in  $[N]$ .

The proof of the above theorem can be found in Appendix D.

In the rest of the paper, we will set  $u = k$  and therefore  $d = 4^k$ . Also we set the smoothness parameter  $J = d^{17} = 4^{17k}$ .

## 5.2 Reduction from Smooth $k$ -Label Cover

The starting point is a smooth  $k$ -LABEL COVER  $\mathcal{L}(G(V, E), M, N, \{\pi^{v, e} | e \in E, v \in e\})$  with  $M = 7^{(J+1)u}$  and  $N = 2^u 7^{Ju}$  as described in Theorem 5.1. Figure 5.2 illustrates the reduction from  $k$ -LABEL COVER  $\mathcal{L}(G(V, E), N, M, \{\pi^{v, e} | e \in E, v \in e\})$  that given an instance of  $k$ -LABEL COVER  $\mathcal{L}$  produces a random labeled example. We refer to the obtained distribution on examples as  $\mathcal{E}$ .

## 5.3 Proof of Theorem 1.1

We claim that our reduction has the following completeness and soundness properties.

**Theorem 5.2.** • **COMPLETENESS:** *If  $\mathcal{L}$  is a strongly-satisfiable instance of smooth  $k$ -LABEL COVER, then there is a disjunction that agrees with a random example from  $\mathcal{E}$  with probability at least  $1 - O(\frac{1}{\sqrt{k}})$ .*

- **SOUNDNESS:** *If  $\mathcal{L}$  is not  $2k^2 2^{-\gamma k}$ -weakly satisfiable and is smooth with parameters  $J = 4^{17k}$  and  $d = 4^k$ , then there is no halfspace that agrees with a random example from  $\mathcal{E}$  with probability more than  $\frac{1}{2} + O(\frac{1}{\sqrt{k}})$ .*

Combining the above theorem with Theorem 5.1 we get that for  $k = O(1/\epsilon^2)$ , we obtain our main result: Theorem 1.1.

It remains to check the correctness of the completeness and soundness claims in Theorem 5.2. First let us prove the completeness property.

*Proof.* (Proof of Completeness) Let  $\Lambda$  be the labeling that strongly satisfies  $\mathcal{L}$ . Consider disjunction  $h(\mathbf{y}) = \bigvee_{v \in V} y_v^{(\Lambda(v))}$ . Let  $e = (v_1, v_2, \dots, v_k)$  be any hyperedge and let  $\mathcal{E}_e$  be the distribution  $\mathcal{E}$  restricted to the examples generated for  $e$ . With probability at least  $1 - 1/k$ ,  $y_{v_i}^{\Lambda(v_i)} = x_i^{\pi^{v_i, e}(\Lambda(v_i))}$  for every  $i \in [k]$ . As  $e$  is strongly satisfied by  $\Lambda$ , for all  $i, j \in [k]$ ,  $\pi^{v_i, e}(\Lambda(v_i)) = \pi^{v_j, e}(\Lambda(v_j))$ . Therefore, as in the proof of Theorem 4.5, we obtain that  $h(\mathbf{y})$  agrees with a random example from  $\mathcal{E}_e$  with probability at least  $1 - O(1/\sqrt{k})$ . Labeling  $\Lambda$  strongly satisfies all edges and therefore we obtain that  $h(\mathbf{y})$  agrees with a random example from  $\mathcal{E}$  with probability at least  $1 - O(1/\sqrt{k})$ .  $\square$

The more complicated part is the soundness property which we prove in Section 5.4.

## 5.4 Soundness Analysis

**Proof Idea** The main idea is similar to the proof of Theorem 4.6 although it is more technically involved. Notice that the reduction in Figure 5.2 produces examples such that  $y_{v_i}^{j_1}, y_{v_i}^{j_2}$  are “almost identical” copies when  $\pi^{v_i, e}(j_1) = \pi^{v_i, e}(j_2)$ . Further for different edges  $e$ , the coordinates of  $\mathbf{y}$  will be grouped in different ways, such that each group will have almost identical copies.

To handle these additional complications, the first step of the proof is to show that almost all the hyperedges in smooth  $k$ -LABEL COVER satisfy a certain “niceness” property. After that we generalize the proofs of Lemma 4.7 and Lemma 4.8 under the weaker assumption that most of the hyperedges are “nice”.

The formal definition of “niceness” and the proof that most of the edges are “nice” appear in Section 5.4.1. The generalization of Lemma 4.7 appears in Section 5.4.2. The generalization of Lemma 4.8 appears in Section 5.4.3. All these results are put together into a proof of Theorem 5.2 in Section 5.4.4.

### 5.4.1 Most of the edges are “nice”

Let  $h(\mathbf{y})$  be a halfspace that agrees with more than  $\frac{1}{2} + \frac{1}{\sqrt{k}}$ -fraction of the examples. Suppose,

$$h(\mathbf{y}) = \text{pos} \left( \sum_{v \in V} \langle \mathbf{w}_v, \mathbf{y}_v \rangle - \theta \right).$$

Let  $\tau = \frac{1}{k^{13}}$  and let

$$\mathbf{s}_v = \text{Truncate}(\mathbf{w}_v, C_\tau(\mathbf{w}_v)), \quad \mathbf{l}_v = \mathbf{w}_v - \mathbf{s}_v.$$

**Definition 5.3.** A vertex  $v \in V$  is said to be  $\beta$ -nice with respect to a hyperedge  $e \in E$  containing it if

$$\sum_{i \in [N]} \left( \sum_{j \in \pi^{-1}(i)} |l_v^{(j)}| \right)^4 \leq \beta \|\mathbf{l}_v\|_2^4,$$

where  $\pi : [M] \rightarrow [N]$  is the projection associated with vertex  $v$  and hyperedge  $e$ . A hyperedge  $e = (v_1, v_2, \dots, v_k)$  is  $\beta$ -nice, if for every  $i \in [k]$ , the vertex  $v_i$  is  $\beta$ -nice with respect to  $e$ .

**Lemma 5.4.** *The fraction of  $2\tau$ -nice hyperedges in  $E$  is at least  $1 - O(1/k)$ .*

*Proof.* By definition, we know that  $\mathbf{l}_v$  is  $\tau$ -regular vector. Denote  $I_v = \{i \mid \frac{(l_v^{(i)})^2}{\|\mathbf{l}_v\|_2^2} \geq \frac{1}{d^8}\}$ . By definition  $|I| \leq d^8$ . Notice there are at most  $d^{16}$  pairs of values in  $I \times I$ . By the smoothness property of the  $k$ -LABEL COVER instance, for any vertex  $v$ , at least  $1 - \frac{d^{16}}{J}$  fraction of the hyperedges incident on  $v$  have the following property: for any  $i, j \in I_v$ ,  $\pi^{v,e}(i) \neq \pi^{v,e}(j)$ . If all the vertices in a hyperedge have this property we call it a *good* hyperedge. By an averaging argument, we know that among all hyperedges at least  $1 - \frac{kd^{16}}{J} = 1 - \frac{k}{4^k} \geq 1 - O(\frac{1}{k})$  fraction is *good*.

We will show all these *good* hyperedges are also  $2\tau$ -nice. For a given *good* hyperedge  $e$ , a vertex  $v \in e$ ,  $\pi = \pi^{v,e}$  and  $i \in [N]$ , there is at most one  $j \in \pi^{-1}(i)$  such that  $\frac{(l_v^{(i)})^2}{\|\mathbf{l}_v\|_2^2} \geq \frac{1}{d^8}$ .

Based on the above property, we will show

$$\sum_{i \in [N]} \left( \sum_{j \in \pi^{-1}(i)} |l_v^{(j)}| \right)^4 \leq 2\tau \|\mathbf{l}_v\|_2^4.$$

Notice that

$$\sum_{i \in [N]} \left( \sum_{j \in \pi^{-1}(i)} |l_v^{(j)}| \right)^4 = \sum_{i \in [N]} \sum_{j_1, j_2, j_3, j_4 \in \pi^{-1}(i)} |l_v^{(j_1)} l_v^{(j_2)} l_v^{(j_3)} l_v^{(j_4)}| \quad (4)$$

and the sum of all the terms with  $j_1 = j_2 = j_3 = j_4$  is  $\|\mathbf{l}_v\|_4^4$ .

For all other terms  $|l_v^{(j_1)} l_v^{(j_2)} l_v^{(j_3)} l_v^{(j_4)}|$  with  $j_1, j_2, j_3, j_4$  that are not all equal, there is at least one  $|l_v^{(j_r)}|$  ( $r \in [4]$ ) smaller than  $\frac{\|\mathbf{l}_v\|_2}{d^4}$ . Therefore,  $|l_v^{(j_1)} l_v^{(j_2)} l_v^{(j_3)} l_v^{(j_4)}|$  can be bounded by

$$\frac{\|\mathbf{l}_v\|_2}{d^4} \left( \sum_{j_1, j_2, j_3, j_4} |l_v^{(j_1)}|^3 + |l_v^{(j_2)}|^3 + |l_v^{(j_3)}|^3 + |l_v^{(j_4)}|^3 \right).$$

Overall, expression (4) can be bounded by

$$\begin{aligned} & \|\mathbf{l}_v\|_4^4 + \frac{\|\mathbf{l}_v\|_2}{d^4} \sum_{i \in [N]} \sum_{j_1, j_2, j_3, j_4 \in \pi^{-1}(i)} |l_v^{(j_1)}|^3 + |l_v^{(j_2)}|^3 + |l_v^{(j_3)}|^3 + |l_v^{(j_4)}|^3 \\ & \leq \tau^2 \|\mathbf{l}_v\|_2^4 + \frac{\|\mathbf{l}_v\|_2}{d^4} 4d^3 \sum_{j \in [M]} |l_v^{(j)}|^3 \quad (\text{since } |\pi^{-1}(i)| \leq d, \text{ each } l_v^{(j)} \text{ appears at most } 4d^3 \text{ times}) \\ & \leq (\tau^2 + 4\frac{\tau}{d}) \|\mathbf{l}_v\|_2^4 \quad (\mathbf{l}_v \text{ is } \tau\text{-regular vector, so } |l_v^j| \leq \tau \|\mathbf{l}_v\|_2 \text{ for all } j \in [M]) \\ & \leq 2\tau \|\mathbf{l}_v\|_2^4. \end{aligned}$$

□

Let us fix a  $2\tau$ -nice hyperedge  $e = (v_1, \dots, v_k)$ . As before let  $\mathcal{E}_e$  denote the distribution on examples restricted to those generated for hyperedge  $e$ . We will analyze the probability that the halfspace  $h(\mathbf{y})$  agrees with a random example from  $\mathcal{E}_e$ .

Let  $\pi^{v_1, e}, \pi^{v_2, e}, \dots, \pi^{v_k, e} : [M] \rightarrow [N]$  denote the projections associated with the hyperedge  $e$ . For the sake of brevity, we shall write  $\mathbf{w}_i, \mathbf{y}_i, \mathbf{l}_i$  instead of  $\mathbf{w}_{v_i}, \mathbf{y}_{v_i}, \mathbf{l}_{v_i}$ . For all  $j \in [N]$  and  $i \in [k]$ , define

$$\mathbf{y}_i^{\{j\}} = \text{Truncate}(\mathbf{y}_i, (\pi^{v_i, e})^{-1}(j)).$$

Similarly, define vectors  $\mathbf{w}_i^{\{j\}}, \mathbf{l}_i^{\{j\}}$  and  $\mathbf{s}_i^{\{j\}}$ .

Notice that for every example  $(\mathbf{y}, b)$  in the support of  $\mathcal{E}_e$ ,  $\mathbf{y}_v = \mathbf{0}$  for every vertex  $v \notin e$ . Therefore, on restricting to examples from  $\mathcal{E}_e$  we can write:

$$h(\mathbf{y}) = \text{pos}\left(\sum_{i \in [k]} \langle \mathbf{w}_i, \mathbf{y}_i \rangle - \theta\right).$$

#### 5.4.2 Common Influential Variables (generalization of Lemma 4.7)

**Lemma 5.5.** *Let  $h(\mathbf{y})$  be a halfspace such that for all  $i \neq j \in [k]$ , we have  $\pi^{v_i, e}(C_\tau(\mathbf{w}_i)) \cap \pi^{v_j, e}(C_\tau(\mathbf{w}_j)) = \emptyset$ . Then*

$$\left| \mathbf{E}_{\mathcal{E}_e}[h(\mathbf{y})|b=0] - \mathbf{E}_{\mathcal{E}_e}[h(\mathbf{y})|b=1] \right| \leq O\left(\frac{1}{k}\right). \quad (5)$$

*Proof.* Fix the following notation:

$$\begin{aligned} \mathbf{y}_i^C &= \text{Truncate}(\mathbf{y}_i, C_\tau(\mathbf{w}_i)) & \mathbf{y}^C &= \mathbf{y}_1^C, \mathbf{y}_2^C, \dots, \mathbf{y}_k^C \\ \mathbf{s} &= \mathbf{s}_1, \mathbf{s}_2, \dots, \mathbf{s}_k & \mathbf{l} &= \mathbf{l}_1, \mathbf{l}_2, \dots, \mathbf{l}_k. \end{aligned}$$

We can rewrite the halfspace  $h(\mathbf{y})$  as  $h(\mathbf{y}) = \text{pos}\left(\langle \mathbf{s}, \mathbf{y}^C \rangle + \langle \mathbf{l}, \mathbf{y} \rangle - \theta\right)$ . Let us first normalize the weights of  $h(\mathbf{y})$  so that  $\sum_{i \in [k]} \|\mathbf{l}_i\|_2^2 = 1$ . Let us condition on a possible fixing of the vector  $\mathbf{y}^C$ . Under this conditioning and also for  $b = 0$ , define the family of ensembles  $\mathcal{A} = \mathbf{A}^{\{1\}}, \dots, \mathbf{A}^{\{N\}}$  as follows:

$$\mathbf{A}^{\{j\}} = \left\{ \mathbf{y}_i^{(r)} \mid i \in [k], r \in [M] \text{ such that } \pi^{v_i, e}(r) = j \text{ and } r \notin C_\tau(\mathbf{w}_i) \right\}$$

Similarly define the ensemble  $\mathcal{B} = \mathbf{B}^{\{1\}}, \dots, \mathbf{B}^{\{N\}}$  for the conditioning  $b = 1$ . Now we shall apply the invariance principle (Theorem 3.10) to the ensembles  $\mathcal{A}, \mathcal{B}$  and the linear function  $\mathbf{l}(\mathbf{y})$ :

$$\mathbf{l}(\mathbf{y}) = \sum_{j \in [N]} \langle \mathbf{l}^{\{j\}}, \mathbf{y}^{\{j\}} \rangle.$$

As we prove in Claim 5.6 below, the ensembles  $\mathcal{A}, \mathcal{B}$  have matching moments up to degree 3. Furthermore, by Lemma 3.3, the linear function  $\mathbf{l}$  and the ensembles  $\mathcal{A}, \mathcal{B}$  satisfy the following spread property:

$$\Pr_{\mathcal{A}} \left[ \mathbf{l}(\mathcal{A}) \in [\theta' - \alpha, \theta' + \alpha] \right] \leq c(\alpha) \quad \Pr_{\mathcal{B}} \left[ \mathbf{l}(\mathcal{B}) \in [\theta' - \alpha, \theta' + \alpha] \right] \leq c(\alpha)$$

for all  $\theta' \in \mathbb{R}$ , where  $c(\alpha) = 8\alpha k + 4\tau k + 2e^{-\frac{1}{2k^4\tau^2}}$  (by setting  $\gamma = \frac{1}{k^2}$  and  $|b - a| = 2\alpha$  in Lemma 3.3).

Using the invariance principle (Th. 3.10), this implies:

$$\left| \mathbf{E}_{\mathcal{A}} \left[ \text{pos} \left( \langle \mathbf{s}, \mathbf{y}^C \rangle + \sum_{j \in [N]} \langle \mathbf{l}^{\{j\}}, \mathbf{A}^{\{j\}} \rangle - \theta \right) | \mathbf{y}^C \right] - \mathbf{E}_{\mathcal{B}} \left[ \text{pos} \left( \langle \mathbf{s}, \mathbf{y}^C \rangle + \sum_{j \in [N]} \langle \mathbf{l}^{\{j\}}, \mathbf{B}^{\{j\}} \rangle - \theta \right) | \mathbf{y}^C \right] \right| \leq O\left(\frac{1}{\alpha^4}\right) \sum_{j \in [N]} \|\mathbf{l}^{\{j\}}\|_1^4 + 2c(\alpha). \quad (6)$$

Take  $\alpha$  to be  $\frac{1}{k^2}$  and recall that  $\tau = \frac{1}{k^{13}}$ . In Claim 5.7 below we show that

$$\sum_{j \in [N]} \|\mathbf{l}^{\{j\}}\|_1^4 \leq 2\tau k^4.$$

The above inequality holds for an arbitrary conditioning of the values of  $\mathbf{y}^C$ . Hence, by averaging over all settings of  $\mathbf{y}^C$  we prove (5).  $\square$

**Claim 5.6.** *The ensembles  $\mathcal{A}$  and  $\mathcal{B}$  have matching moments up to degree 3.*

Let us suppose for a moment that  $\mathbf{y}$  was generated by setting  $y_{v_i}^{(j)} = x_i^{(\pi^{v_i, e(j)})}$ , that is without adding any noise. By Lemma 4.1, the first four moments of random variable  $\mathbf{y}$  conditioned on  $b = 0$  agree with the first moments of random variable  $\mathbf{y}$  conditioned on  $b = 1$ . As we showed in Observation 4.3, even with noise, the first four moments of  $\mathbf{y}$  remain the same when conditioned on  $b = 0$  and  $b = 1$ . Finally,  $\pi^{v_i, e}(C_\tau(\mathbf{w}_i)) \cap \pi^{v_j, e}(C_\tau(\mathbf{w}_j)) = \emptyset$  for all  $i \neq j \in [k]$ . Hence for each  $j \in [N]$ , conditioning on  $\mathbf{y}^C$  fixes bits in at most one row of  $\mathbf{A}^{\{j\}}$ . Formally, for every  $j \in [N]$ , there exists at most one  $i \in [k]$  such that  $\mathbf{y}_i^{\{j\}}$  and  $\mathbf{y}^C$  have shared variables. Therefore, by Lemma 4.4,  $\mathcal{A}$  and  $\mathcal{B}$  have matching moments up to degree 3.

**Claim 5.7.**

$$\sum_{j \in [N]} \|\mathbf{l}^{\{j\}}\|_1^4 \leq 2\tau k^4.$$

*Proof.* Since  $\|\mathbf{l}^{\{j\}}\|_1 = \sum_{i \in [k]} \|\mathbf{l}_i^{\{j\}}\|_1$ , we can write

$$\sum_{j \in [N]} \|\mathbf{l}^{\{j\}}\|_1^4 \leq \sum_{j \in [N]} k^4 \left( \sum_{i \in [k]} \|\mathbf{l}_i^{\{j\}}\|_1^4 \right) = k^4 \sum_{i \in [k]} \left( \sum_{j \in [N]} \|\mathbf{l}_i^{\{j\}}\|_1^4 \right). \quad (7)$$

As  $e = (v_1, \dots, v_k)$  is a  $2\tau$ -nice hyperedge, we have  $\sum_{j \in [N]} \|\mathbf{l}_i^{\{j\}}\|_1^4 \leq 2\tau \|\mathbf{l}_i\|_2^4$ . By normalization of  $\mathbf{l}$ , we know  $\sum_{i \in [k]} \|\mathbf{l}_i\|_2^2 = 1$ . Substituting this into inequality (7) we get the claimed bound.  $\square$

### 5.4.3 Bounding the Number of Influential Coordinates (generalization of Lemma 4.8)

**Lemma 5.8.** *Given a halfspace  $h(\mathbf{y}) = \text{pos}(\sum_{i \in [k]} \langle \mathbf{w}_i, \mathbf{y}_i \rangle - \theta)$  and  $r \in [k]$  such that  $|C_\tau(\mathbf{w}_r)| \geq t$  for  $t = \frac{1}{\tau^2} (\lceil 4k^2 \ln(2k) \rceil \lceil 4 \ln(1/\tau) \rceil + \ln(1/\tau) + 10 \ln d) = O(k^{29})$ , define  $\tilde{h}(\mathbf{y}) = \text{pos}(\sum_{i \in [k]} \langle \tilde{\mathbf{w}}_i, \mathbf{y}_i \rangle - \tilde{\theta})$  as follows:*

- $\tilde{\mathbf{w}}_r = \text{Truncate}(\mathbf{w}_r, H_t(\mathbf{w}_r))$  and  $\tilde{\mathbf{w}}_i = \mathbf{w}_i$  for all  $i \neq r$ .
- $\tilde{\theta} = \theta - \mathbf{E}[\langle \mathbf{a}_r, \mathbf{y}_r \rangle | b = 0]$ , for  $\mathbf{a} = \mathbf{w} - \tilde{\mathbf{w}}$ .

Then,

$$\left| \mathbf{E}_{\mathcal{E}_e}[\tilde{h}(\mathbf{y}) | b = 0] - \mathbf{E}_{\mathcal{E}_e}[h(\mathbf{y}) | b = 0] \right| \leq \frac{1}{k^2}, \quad \left| \mathbf{E}_{\mathcal{E}_e}[\tilde{h}(\mathbf{y}) | b = 1] - \mathbf{E}_{\mathcal{E}_e}[h(\mathbf{y}) | b = 1] \right| \leq \frac{1}{k^2}.$$

*Proof.* It is easy to see that the matching moments condition implies that

$$\mathbf{E}_{\mathcal{E}_e}[\langle \mathbf{a}_r, \mathbf{y}_r \rangle | b = 0] = \mathbf{E}_{\mathcal{E}_e}[\langle \mathbf{a}_r, \mathbf{y}_r \rangle | b = 1].$$

Let us show the inequality for the case  $b = 0$ , the other inequality can be derived in an identical way. Let  $\mathcal{E}_{e,0}$  denote distribution  $\mathcal{E}_e$  conditioned on  $b = 0$ . Without loss of generality, we may assume that  $r = 1$  and  $|w_1^{(1)}| \geq |w_1^{(2)}| \dots \geq |w_1^{(M)}|$ . In particular, this implies  $H_t(\mathbf{w}_1) = \{1, \dots, t\}$ . Define

$$\mu_r = \mathbf{E}_{\mathcal{E}_{e,0}}[\langle \mathbf{a}_r, \mathbf{y}_r \rangle], \quad \mu_r^{\{i\}} = \mathbf{E}_{\mathcal{E}_{e,0}}[\langle \mathbf{a}_r^{\{i\}}, \mathbf{y}_r^{\{i\}} \rangle].$$

Let us set  $T = \lceil 4k^2 \ln(2k) \rceil$  and define the subset  $G = \{g_1, \dots, g_T\}$  of  $H_t(\mathbf{w}_1)$  as follows:

$$G = \{g_i \mid g_i = 1 + i \lceil (4/\tau^2) \ln(1/\tau) \rceil, 0 \leq i \leq T\}.$$

Therefore, by Lemma 3.2,  $|w_1^{(g_i)}|$  is a geometrically decreasing sequence such that  $|w_1^{(g_{i+1})}| \leq |w_1^{(g_i)}|/3$ . Let  $H = H_t(\mathbf{w}_1) \setminus G$ . Fix the following notation:

$$\mathbf{w}_1^G = \text{Truncate}(\mathbf{w}_1, G), \quad \mathbf{w}_1^H = \text{Truncate}(\mathbf{w}_1, H), \quad \mathbf{w}_1^{>t} = \text{Truncate}(\mathbf{w}_1, \{t+1, \dots, n\}).$$

Similarly, define the vectors  $\mathbf{y}_1^G, \mathbf{y}_1^H, \mathbf{y}_1^{>t}$ . By definition, we have  $\mathbf{a}_1 = \mathbf{w}_1^{>t}$ . Rewriting the halfspace functions  $h(\mathbf{y}), \tilde{h}(\mathbf{y})$ :

$$\begin{aligned} h(\mathbf{y}) &= \text{pos} \left( \sum_{i=2}^k \langle \mathbf{w}_i, \mathbf{y}_i \rangle + \langle \mathbf{w}_1^G, \mathbf{y}_1^G \rangle + \langle \mathbf{w}_1^H, \mathbf{y}_1^H \rangle + \langle \mathbf{a}_1, \mathbf{y}_1^{>t} \rangle - \theta \right), \\ \tilde{h}(\mathbf{y}) &= \text{pos} \left( \sum_{i=2}^k \langle \mathbf{w}_i, \mathbf{y}_i \rangle + \langle \mathbf{w}_1^G, \mathbf{y}_1^G \rangle + \langle \mathbf{w}_1^H, \mathbf{y}_1^H \rangle + \mu_1 - \theta \right). \end{aligned}$$

By Claim 5.9 below, with probability at most  $\frac{1}{d} = \frac{1}{4^k}$ , we have  $|\langle \mathbf{a}_1, \mathbf{y}_1 \rangle - \mu_1| \geq d^4 \|\mathbf{a}_1\|_2$ . Suppose  $|\langle \mathbf{a}_1, \mathbf{y}_1 \rangle - \mu_1| < d^4 \|\mathbf{a}_1\|_2$ , then Claim 5.10 below gives  $|\langle \mathbf{a}_1, \mathbf{y}_1 \rangle - \mu_1| < 1/d^6 |w_1^{(gt)}| < \frac{1}{3} |w_1^{(gt)}|$ . Thus, we can write

$$\Pr_{\mathcal{E}_{e,0}} [h(\mathbf{y}) \neq \tilde{h}(\mathbf{y})] \leq \Pr_{\mathcal{E}_{e,0}} \left[ \langle \mathbf{w}_1^G, \mathbf{y}_1^G \rangle \in [\theta' - \frac{1}{3} |w_1^{(gt)}|, \theta' + \frac{1}{3} |w_1^{(gt)}|] \right] + \frac{1}{4^k}.$$

where  $\theta' = -\sum_{i=2}^k \langle \mathbf{w}_i, \mathbf{y}_i \rangle - \langle \mathbf{w}_1^H, \mathbf{y}_1^H \rangle - \mu_1 + \theta$ . For any fixing of the value of  $\theta' \in \mathbb{R}$ , induces a certain distribution on  $\mathbf{y}_1^G$ . However, the  $\frac{1}{k^2}$  noise introduced in  $\mathbf{y}_1^G$  is completely independent. This corresponds to the setting of Lemma 3.7, and hence we can bound the above probability by  $(1 - 1/(2k^2))^T + 1/4^k \leq (1 - 1/(2k^2))^{4k^2 \ln(2k)} + 1/4^k \leq 1/k^2$ .  $\square$

**Claim 5.9.**

$$\Pr_{\mathcal{E}_{e,0}} \left[ |\langle \mathbf{a}_1, \mathbf{y}_1 \rangle - \mu_1| \geq d^4 \|\mathbf{a}_1\|_2 \right] \leq \frac{1}{d}.$$

*Proof.* Write  $[M]$  as the union of disjoint sets  $R_1 \cup R_2 \cup \dots \cup R_N$  where  $R_i = (\pi^{v_i, e})^{-1}(i)$ . Notice every  $R_i$  has size at most  $d$ , therefore

$$\mathbf{Var}_{\mathcal{E}_{e,0}}(\langle \mathbf{a}_1, \mathbf{y}_1 \rangle) = \sum_{i \in [N]} \mathbf{Var}_{\mathcal{E}_{e,0}}(\langle \mathbf{a}_1^{R_i}, \mathbf{y}_1^{R_i} \rangle) \leq \sum_{i \in [N]} d \|\mathbf{a}_1^{R_i}\|_2^2 = d \|\mathbf{a}_1\|_2^2.$$

By applying Chebyshev's inequality (Th. A.3), we have

$$\Pr_{\mathcal{E}_{e,0}} \left[ |\langle \mathbf{a}_1, \mathbf{y}_1 \rangle - \mu_1| \geq d^4 \|\mathbf{a}_1\|_2 \right] \leq \frac{2d}{d^8} \leq \frac{1}{d}.$$

□

**Claim 5.10.** *By the choice of the parameters  $T$  and  $t$ ,*

$$\|\mathbf{a}_1\|_2 \leq \frac{1}{d^{10}} |w_1^{(gT)}|.$$

*Proof.* By Lemma 3.2,

$$|w_1^{(gT)}|^2 \geq \frac{\tau}{(1-\tau^2)^{t-gT}} \|\mathbf{a}_1\|_2^2 \geq \frac{\tau}{(1-\tau^2)^{\frac{1}{\tau^2}(\ln(1/\tau)+10 \ln d)}} \|\mathbf{a}_1\|_2^2 \geq d^{10} \|\mathbf{a}_1\|_2^2.$$

□

#### 5.4.4 Proof of Soundness

Recall that we chose  $\tau = 1/k^{13}$  and  $t = O(k^{29})$ .

**Lemma 5.11.** *Fix a hyperedge  $e$  which is  $2\tau$ -nice. If for all  $i \neq j \in [k]$ ,  $\pi^{v_i, e}(H_t(\mathbf{w}_i)) \cap \pi^{v_j, e}(H_t(\mathbf{w}_j)) = \emptyset$  then the probability that halfspace  $h(\mathbf{y})$  agrees with a random example from  $\mathcal{E}_e$  is at most  $\frac{1}{2} + O(\frac{1}{k})$ .*

*Proof.* The proof is similar to the proof of Theorem 4.6. Define  $I = \{r \mid C_\tau(\mathbf{w}_r) > t\}$ . We divide the problem into the following two cases.

1.  $I = \emptyset$ ; i.e., for all  $i \in [k]$ ,  $C_\tau(\mathbf{w}_i) \leq t$ . Then for any  $i \neq j \in [k]$ ,  $H_t(\mathbf{w}_i) \cap H_t(\mathbf{w}_j) = \emptyset$  implies  $C_\tau(\mathbf{w}_i) \cap C_\tau(\mathbf{w}_j) = \emptyset$ . By Lemma 5.5, we have

$$\left| \mathbf{E}_{\mathcal{E}_e}[h(\mathbf{y})|b=0] - \mathbf{E}_{\mathcal{E}_e}[h(\mathbf{y})|b=1] \right| \leq O\left(\frac{1}{k}\right).$$

2.  $I \neq \emptyset$ . Then for all  $r \in I$ , we set  $\tilde{\mathbf{w}}_r = \text{Truncate}(\mathbf{w}_r, H_t(\mathbf{w}_r))$  and define a new halfspace  $h'$  by replacing  $\mathbf{w}_r$  with  $\tilde{\mathbf{w}}_r$  in  $h$ . Since such replacements occur at most  $k$  times and, by Lemma

5.8, every replacement changes the output of the halfspace on at most  $\frac{1}{k^2}$  fraction of examples from  $\mathcal{E}_e$ , we can bound the overall change by  $k \times \frac{1}{k^2} = \frac{1}{k}$ . That is

$$\left| \mathbf{E}_{\mathcal{E}_{e,0}} [h'(\mathbf{y})] - \mathbf{E}_{\mathcal{E}_{e,0}} [h(\mathbf{y})] \right| \leq \frac{1}{k}, \quad \left| \mathbf{E}_{\mathcal{E}_{e,1}} [h'(\mathbf{y})] - \mathbf{E}_{\mathcal{E}_{e,1}} [h(\mathbf{y})] \right| \leq \frac{1}{k}. \quad (8)$$

For the halfspace  $h'$  and for all  $r \in [k]$ , we have  $|C_\tau(\tilde{\mathbf{w}}_r)| \leq t$ , thus reducing to Case 1. Therefore

$$\left| \mathbf{E}_{\mathcal{E}_{e,o}} [h'(\mathbf{y})] - \mathbf{E}_{\mathcal{E}_{e,1}} [h'(\mathbf{y})] \right| \leq O\left(\frac{1}{k}\right). \quad (9)$$

Combining (8) and (9), we get

$$\left| \mathbf{E}_{\mathcal{E}_{e,0}} [h(\mathbf{y})] - \mathbf{E}_{\mathcal{E}_{e,1}} [h(\mathbf{y})] \right| \leq O\left(\frac{1}{k}\right).$$

In other words, the probability that halfspace  $h(\mathbf{y})$  agrees with a random example from  $\mathcal{E}_e$  is at most  $\frac{1}{2} + O\left(\frac{1}{k}\right)$ .  $\square$

We first recall the soundness statement:

**Proposition 5.12.** *If  $\mathcal{L}$  is not a  $2k^2 2^{-\gamma k}$ -weakly satisfiable instance of smooth  $k$ -LABEL COVER, then there is no halfspace that agrees with a random example from  $\mathcal{E}$  with probability more than  $\frac{1}{2} + \frac{1}{\sqrt{k}}$ .*

*Proof.* The proof is by contradiction. We can define the following labeling strategy: for each vertex  $v$ , uniformly randomly pick a label from  $H_t(\mathbf{w}_v)$ . We know that the size of  $H_t(\mathbf{w}_{v_i})$  is  $t = O(k^{29})$ .

Suppose there exists a halfspace that agrees with a random example from  $\mathcal{E}$  with probability more than  $\frac{1}{2} + \frac{1}{\sqrt{k}}$ . Then by an averaging argument, for at least  $\frac{1}{2\sqrt{k}}$ -fraction of the hyperedges  $e$ ,  $h(\mathbf{y})$  agrees with a random example from  $\mathcal{E}_e$  with probability at least  $\frac{1}{2} + \frac{1}{2\sqrt{k}}$ . We refer to these edges as *good*.

Since there is at most  $O(1/k)$ -fraction of the hyperedges that are *not*  $2\tau$ -nice we know that at least  $\frac{1}{4\sqrt{k}}$ -fraction of the hyperedges are  $2\tau$ -nice and *good*. By Lemma 5.11, for each  $2\tau$ -nice and *good* hyperedge  $e$  there exist two vertices  $v_i, v_j \in e$  such that  $\pi^{v_i, e}(H_t(\mathbf{w}_i))$  and  $\pi^{v_j, e}(H_t(\mathbf{w}_j))$  intersect. Then there is a  $\frac{1}{t^2}$  probability that the labeling strategy we defined will weakly satisfy hyperedge  $e$ .

Overall this strategy is expected to weakly satisfy at least  $\frac{1}{4\sqrt{k}} \frac{1}{t^2} = \Omega\left(\frac{1}{k^{59}}\right)$  fraction of the hyperedges. This is a contradiction since  $\mathcal{L}$  is not  $\frac{2k^2}{2^{\gamma k}}$ -weakly satisfiable.  $\square$

## References

- [1] E. Amaldi and V. Kann. On the approximability of minimizing nonzero variables or unsatisfied relations in linear systems. *Theoretical Computer Science*, 109:237–260, 1998. [3](#)
- [2] D. Angluin and P. Laird. Learning from noisy examples. *Machine Learning*, 2:343–370, 1988. [3](#)

- [3] S. Arora, L. Babai, J. Stern, and Z. Sweedyk. The hardness of approximate optima in lattices, codes, and systems of linear equations. *J. Comput. Syst. Sci.*, 54(2):317–331, 1997. [3](#)
- [4] P. Auer and M. K. Warmuth. Tracking the best disjunction. *Machine Learning*, 32(2):127–150, 1998. [4](#)
- [5] S. Ben-David, N. Eiron, and P. M. Long. On the difficulty of approximately maximizing agreements. *J. Comput. Syst. Sci.*, 66(3):496–514, 2003. [3](#)
- [6] A. Blum, A. Frieze, R. Kannan, and S. Vempala. A polynomial-time algorithm for learning noisy linear threshold functions. *Algorithmica*, 22(1-2):35–52, 1998. [3](#)
- [7] A. Blumer, A. Ehrenfeucht, D. Haussler, and M. K. Warmuth. Occam’s razor. *Inf. Process. Lett.*, 24(6):377–380, 1987. [2](#)
- [8] N. Bshouty and L. Burroughs. Bounds for the minimum disagreement problem with applications to learning theory. In *Proceedings of COLT*, pages 271–286, 2002. [3](#)
- [9] N. Bshouty and L. Burroughs. Maximizing agreements and coagnostic learning. *Theoretical Computer Science*, 350(1):24–39, 2006. [3](#)
- [10] T. Bylander. Learning linear threshold functions in the presence of classification noise. In *Proceedings of COLT*, pages 340–347, 1994. [3](#)
- [11] S. Chatterjee. A simple invariance theorem. *arxiv:math/0508213v1.*, 2005. [11](#)
- [12] E. Cohen. Learning noisy perceptrons by a perceptron in polynomial time. In *IEEE FOCS*, pages 514–523, 1997. [3](#)
- [13] I. Diakonikolas, P. Gopalan, R. Jaiswal, R. A. Servedio, and E. Viola. Bounded independence fools halfspaces. In *FOCS*, pages 171–180, 2009. [6](#), [7](#), [9](#), [10](#)
- [14] V. Feldman. Optimal hardness results for maximizing agreements with monomials. In *IEEE CCC*, pages 226–236, 2006. [3](#)
- [15] V. Feldman, P. Gopalan, S. Khot, and A. K. Ponnuswami. On agnostic learning of parities, monomials, and halfspaces. *SIAM J. Comput.*, 39(2):606–645, 2009. [3](#)
- [16] S. Galant. Perceptron based learning algorithms. *IEEE Trans. on Neural Networks*, 1(2), 1990. [4](#)
- [17] M. Garey and D. S. Johnson. *Computers and Intractability*. 1979. [3](#)
- [18] C. Gentile and M. K. Warmuth. Linear hinge loss and average margin. In *Proceedings of NIPS*, pages 225–231, 1998. [4](#)
- [19] P. Gopalan, S. Khot, and R. Saket. Hardness of reconstructing multivariate polynomials over finite fields. *SIAM J. Comput.*, 39(6):2598–2621, 2010. [36](#)
- [20] V. Guruswami and P. Raghavendra. Hardness of learning halfspaces with noise. *SIAM J. Comput.*, 39(2):742–765, 2009. [3](#)

- [21] D. Haussler. Decision theoretic generalizations of the PAC model for neural net and other learning applications. *Information and Computation*, 100(1):78–150, 1992. [2](#)
- [22] K. Hoffgen, K. van Horn, and H. U. Simon. Robust trainability of single neurons. *J. Comput. Syst. Sci.*, 50(1):114–125, 1995. [3](#)
- [23] D. S. Johnson and F. P. Preparata. The densest hemisphere problem. *Theoretical Computer Science*, 6:93–107, 1978. [3](#)
- [24] A. Kalai, A. Klivans, Y. Mansour, and R. Servedio. Agnostically learning halfspaces. *SIAM Journal on Computing*, 37(6):1777–1805, 2008. [3](#), [4](#)
- [25] M. Kearns. Efficient noise-tolerant learning from statistical queries. *Journal of the ACM*, 45(6):983–1006, 1998. [3](#)
- [26] M. Kearns, R. Schapire, and L. Sellie. Toward efficient agnostic learning. *Machine Learning*, 17:115–141, 1994. [2](#), [3](#)
- [27] M. J. Kearns and M. Li. Learning in the presence of malicious errors. *SIAM J. Comput.*, 22(4):807–837, 1993. [3](#)
- [28] M. J. Kearns and R. E. Schapire. Efficient distribution-free learning of probabilistic concepts. *J. Comput. Syst. Sci.*, 48(3):464–497, 1994. [3](#)
- [29] S. Khot. On the power of unique 2-Prover 1-Round games. In *ACM STOC*, pages 767–775, May 19–21 2002. [4](#), [5](#)
- [30] S. Khot. New techniques for probabilistically checkable proofs and inapproximability results (thesis). *Princeton University Technical Reports*, TR-673-03, 2003. [8](#), [36](#)
- [31] S. Khot, G. Kindler, E. Mossel, and R. O’Donnell. Optimal inapproximability results for MAX-CUT and other 2-variable CSPs? *SIAM J. Comput.*, 37(1):319–357, 2007. [8](#)
- [32] S. Khot and R. Saket. Hardness of minimizing and learning DNF expressions. In *IEEE FOCS*, pages 231–240, 2008. [2](#), [3](#)
- [33] W. S. Lee, P. L. Bartlett, and R. C. Williamson. On efficient agnostic learning of linear combinations of basis functions. In *Proceedings of COLT*, pages 369–376, 1995. [3](#)
- [34] N. Littlestone. Learning quickly when irrelevant attributes abound: A new linear-threshold algorithm. *Machine Learning*, 2:285–318, 1987. [2](#)
- [35] K. Matulef, R. O’Donnell, R. Rubinfeld, and R. A. Servedio. Testing halfspaces. In *SODA*, pages 256–264, 2009. [6](#), [9](#)
- [36] E. Mossel. Gaussian bounds for noise correlation of functions. *IEEE FOCS*, 2008. [11](#)
- [37] E. Mossel, R. O’Donnell, and K. Oleszkiewicz. Noise stability of functions with low influences: Invariance and optimality. In *IEEE FOCS*, 2005. [11](#), [35](#)
- [38] R. O’Donnell and R. A. Servedio. The chow parameters problem. In *ACM STOC*, pages 517–526, 2008. [6](#), [9](#)

- [39] R. Rivest. Learning decision lists. *Machine Learning*, 2(3):229–246, 1987. [2](#)
- [40] F. Rosenblatt. The perceptron: a probabilistic model for information storage and organization in the brain. *Psychological Review*, 65:386–407, 1958. [2](#)
- [41] R. A. Servedio. Every linear threshold function has a low-weight approximator. *Comput. Complex.*, 16(2):180–209, 2007. [6](#), [9](#)
- [42] L. Valiant. A theory of the learnable. *Communications of the ACM*, 27(11):1134–1142, 1984. [2](#)
- [43] V. Vapnik. *Statistical Learning Theory*. 1998. [2](#)

## Appendix

### A Probabilistic Inequalities

In the discussion below we will make use of the following well-known inequalities.

**Theorem A.1.** (*Hoeffding’s Inequality*) Let  $x^{(1)}, \dots, x^{(n)}$  be independent real random variables such that  $x^{(i)} \in [a^{(i)}, b^{(i)}]$ . Then the sum of these variables  $S = \sum_{i=1}^n x^{(i)}$  satisfies

$$\Pr[|S - \mathbf{E}[S]| \geq nt] \leq 2e^{-\frac{n^2 t^2}{\sum_{i=1}^n (b^{(i)} - a^{(i)})^2}}.$$

**Theorem A.2.** (*Berry-Esseen Theorem*) Let  $x_1, x_2, \dots, x_n$  be i.i.d. random unbiased  $\{-1, 1\}$  variables. Also assume that  $\sum_{i=1}^n c_i^2 = 1$  and  $\max_i \{|c_i|\} \leq \alpha$ . Let  $g$  denote a unit Gaussian variable  $N(0, 1)$ . Then for any  $t \in \mathbb{R}$ ,

$$\left| \Pr \left[ \sum c_i x_i \leq t \right] - \Pr[g \leq t] \right| \leq \alpha.$$

**Theorem A.3.** (*Chebyshev’s Inequality*) Let  $X$  be a random variable with expected value  $\mu$  and variance  $\sigma^2$ . Then for any real number  $t > 0$ ,

$$\Pr[|X - \mu| \geq t \cdot \sigma] \leq 1/t^2.$$

### B Proof of Lemma 3.3

Recall that each  $y^{(i)}$  is generated by the following manner:

$$y^{(i)} = \begin{cases} x^{(i)} & \text{with probability } 1 - \gamma \\ \text{random bit} & \text{with probability } \gamma. \end{cases} \quad (10)$$

Let us define a random vector  $\mathbf{z} \in \{0, 1\}^n$  based on  $\mathbf{y}$ . For  $\mathbf{y}$  generated, if  $y^{(i)}$  is generated as a copy of  $x^{(i)}$  in (10), then  $z^{(i)} = 0$ ; if  $y^{(i)}$  is generated as a random bit in (10), then  $z^{(i)} = 1$ . Let us write  $S = \sum_{i=1}^n w^{(i)} y^{(i)}$ . Our proof is based on two claims.

**Claim B.1.** For a  $\tau$ -regular vector  $\mathbf{w}$ ,  $\Pr[\sum_{i=1}^n |w^{(i)}|^2 z^{(i)} \geq \gamma/2] \geq 1 - 2e^{-\frac{\gamma^2}{2\tau^2}}$ .

**Claim B.2.** For a  $\tau$ -regular vector  $\mathbf{w}$ , given any  $a' < b' \in \mathbb{R}$  and any fixing of  $z^{(1)}, z^{(2)}, \dots, z^{(n)}$ , if  $\sum_{i=1}^n (w^{(i)})^2 z^{(i)} = \sigma^2 > 0$ , then  $\Pr[S \in [a', b']] \leq \frac{2|b'-a'|}{\sigma} + \frac{2\tau}{\sigma}$ .

Given the above two claims are correct, define event  $V$  to be  $\{\sum_{i=1}^n (w^{(i)})^2 z^{(i)} \geq \frac{\gamma}{2}\}$  and use  $\mathbf{1}_{[a,b]}(x) : \mathbb{R} \rightarrow \{0, 1\}$  to denote the indicator function of whether  $x$  falls into interval  $[a, b]$ .

$$\Pr[S \in [a, b]] = \mathbf{E}[\mathbf{1}_{[a,b]}(S)] = \Pr[V] \mathbf{E}[\mathbf{1}_{[a,b]}(S) | V] + \Pr[\neg V] \mathbf{E}[\mathbf{1}_{[a,b]}(S) | \neg V]$$

By Claim B.1,

$$\Pr[\neg V] \mathbf{E}[\mathbf{1}_{[a,b]}(S) | \neg V] \leq \Pr[\neg V] \leq 2e^{-\frac{\gamma^2}{2\tau^2}}.$$

By Claim B.2,

$$\Pr[V] \mathbf{E}[\mathbf{1}_{[a,b]}(S) | V] \leq \frac{4(b-a)}{\sqrt{\gamma}} + \frac{4\tau}{\sqrt{\gamma}}.$$

Overall,

$$\Pr[S \in [a, b]] \leq \frac{4(b-a)}{\sqrt{\gamma}} + \frac{4\tau}{\sqrt{\gamma}} + 2e^{-\frac{\gamma^2}{2\tau^2}}.$$

It remains to verify Claim B.1 and Claim B.2.

To prove Claim B.1, we need to apply the Hoeffding's inequality (see Theorem A.1).

Notice that  $(w^{(i)})^2 z^{(i)} \in [0, (w^{(i)})^2]$  and applying Hoeffding's Inequality, we know

$$\Pr \left[ \left| \sum_{i=1}^n (w^{(i)})^2 z^{(i)} - \mathbf{E} \left[ \sum_{i=1}^n (w^{(i)})^2 z^{(i)} \right] \right| \geq nt \right] \leq 2e^{\frac{-2n^2 t^2}{\sum_{i=1}^n (w^{(i)})^4}}.$$

We know  $\mathbf{E}[\sum_{i=1}^n (w^{(i)})^2 z^{(i)}] = \gamma$  and  $\sum_{i=1}^n ((w^{(i)})^2)^2 \leq \max_i \{(w^{(i)})^2\} \sum_{i=1}^n (w^{(i)})^2 \leq \tau^2$ . If we take  $nt = \gamma/2$ , we have

$$\Pr \left[ \left| \sum_{i=1}^n (w^{(i)})^2 z^{(i)} - \gamma \right| \geq \frac{\gamma}{2} \right] \leq 2e^{-\frac{\gamma^2}{2\tau^2}}.$$

Therefore, with probability at least  $1 - 2e^{-\frac{\gamma^2}{2\tau^2}}$ ,  $\sum_{i=1}^n (w^{(i)})^2 z^{(i)} \geq \frac{\gamma}{2}$ .

To prove Claim B.2, we need use Berry-Esseen Theorem (See Theorem A.2). Let us split  $S$  into two parts:  $S' = \sum_{z_i=1} w_i y_i$  and  $S'' = \sum_{z_i=0} w_i y_i$ . Since  $S = S' + S''$  and  $S'$  is independent of  $S''$ , it suffices to show that  $\Pr[S' \in [a', b']] \leq \frac{2|b'-a'|}{\sqrt{\sigma}} + \frac{2\tau}{\sigma}$  for any  $a', b' \in \mathbb{R}$ . Define  $y'^{(i)} = 2y^{(i)} - 1$  and note that  $y'^{(i)}$  a  $\{-1, 1\}$  variable. By rewriting  $S'$  using this definition, we have

$$S' = \sum_{z^{(i)}=1} w^{(i)} y^{(i)} = \sum_{z^{(i)}=1} w^{(i)} \frac{1 + y'^{(i)}}{2}.$$

Then

$$\Pr [S' \in [a', b']] = \Pr \left[ \sum_{z^{(i)}=1} w^{(i)} y'^{(i)} \in [a'', b''] \right], \quad (11)$$

where  $a'' = 2a' - \sum_{z^{(i)}=1} w^{(i)}$  and  $b'' = 2b' - \sum_{z^{(i)}=1} w^{(i)}$ . We can further rewrite the above term as

$$\begin{aligned} & \Pr \left[ \sum_{z^{(i)}=1} w^{(i)} y'^{(i)} \leq b'' \right] - \Pr \left[ \sum_{z^{(i)}=1} w^{(i)} y'^{(i)} \leq a'' \right] \\ &= \Pr \left[ \sum_{z^{(i)}=1} \frac{w^{(i)} y'^{(i)}}{\sqrt{\sum_{z^{(i)}=1} (w^{(i)})^2}} \leq \frac{b''}{\sqrt{\sum_{z^{(i)}=1} (w^{(i)})^2}} \right] - \Pr \left[ \sum_{z^{(i)}=1} \frac{w^{(i)} y'^{(i)}}{\sqrt{\sum_{z^{(i)}=1} (w^{(i)})^2}} \leq \frac{a''}{\sqrt{\sum_{z^{(i)}=1} (w^{(i)})^2}} \right]. \end{aligned}$$

We can now apply Berry-Esseen's theorem. Notice that for all the  $i$  such that  $z^{(i)} = 1$ ,  $y'^{(i)}$  is distributed as an independent unbiased random  $\{-1, 1\}$  variable. Also  $\max_{z^{(i)}=1} \frac{|w^{(i)}|}{\sqrt{\sum_{z^{(i)}=1} (w^{(i)})^2}} \leq \frac{\tau}{\sqrt{\sum_{z^{(i)}=1} (w^{(i)})^2}}$ .

By Berry-Esseen's theorem, we know that expression (11) is bounded by

$$\Pr \left[ N(0, 1) \leq \frac{b''}{\sqrt{\sum_{z^{(i)}=1} (w^{(i)})^2}} \right] - \Pr \left[ N(0, 1) \leq \frac{a''}{\sqrt{\sum_{z^{(i)}=1} (w^{(i)})^2}} \right] + \frac{2\tau}{\sqrt{\sum_{z^{(i)}=1} (w^{(i)})^2}}.$$

Using the fact that a unit Gaussian variable falls in any interval of length  $\lambda$  with probability at most  $\lambda$  and noticing that  $b'' - a'' = 2(b' - a')$ , we can bound the above quantity by

$$\frac{2|b' - a'|}{\sqrt{\sum_{z^{(i)}=1} (w^{(i)})^2}} + \frac{2\tau}{\sqrt{\sum_{z^{(i)}=1} (w^{(i)})^2}} = \frac{2|b - a|}{\sigma} + \frac{2\tau}{\sigma}.$$

## C Proof of Invariance Principle (Th. 3.10)

We restate our version of the invariance principle here for convenience.

**Theorem 3.10 restated** (Invariance Principle) Let  $\mathcal{A} = \{\mathbf{A}^{\{1\}}, \dots, \mathbf{A}^{\{R\}}\}$ ,  $\mathcal{B} = \{\mathbf{B}^{\{1\}}, \dots, \mathbf{B}^{\{R\}}\}$  be families of ensembles of random variables with  $\mathbf{A}^{\{i\}} = \{a_1^{(i)}, \dots, a_{k_i}^{(i)}\}$  and  $\mathbf{B}^{\{i\}} = \{b_1^{(i)}, \dots, b_{k_i}^{(i)}\}$ , satisfying the following properties:

- For each  $i \in [R]$ , the random variables in ensembles  $(\mathbf{A}^{\{i\}}, \mathbf{B}^{\{i\}})$  have matching moments up to degree 3. Further all the random variables in  $\mathcal{A}$  and  $\mathcal{B}$  are bounded by 1.
- The ensembles  $\mathbf{A}^{\{i\}}$  are all independent of each other, similarly the ensembles  $\mathbf{B}^{\{i\}}$  are independent of each other.

Given a set of vectors  $\mathbf{l} = \{\mathbf{l}^{\{1\}}, \dots, \mathbf{l}^{\{R\}}\}$  ( $\mathbf{l}^{\{i\}} \in \mathbb{R}^{k_i}$ ), define the linear function  $\mathbf{l} : \mathbb{R}^{k_1} \times \dots \times \mathbb{R}^{k_R} \rightarrow \mathbb{R}$  as

$$\mathbf{l}(\mathbf{x}) = \sum_{i \in [R]} \langle \mathbf{l}^{\{i\}}, \mathbf{x}^{\{i\}} \rangle$$

Then for a  $K$ -bounded function  $\Psi : \mathbb{R} \rightarrow \mathbb{R}$  we have

$$\left| \mathbf{E}_{\mathcal{A}} [\Psi(\mathbf{l}(\mathcal{A}) - \theta)] - \mathbf{E}_{\mathcal{B}} [\Psi(\mathbf{l}(\mathcal{B}) - \theta)] \right| \leq K \sum_{i \in [R]} \|\mathbf{l}^{\{i\}}\|_1^4. \quad (12)$$

for all  $\theta > 0$ . Further, define the spread function  $c(\alpha)$  corresponding to the ensembles  $\mathcal{A}, \mathcal{B}$  and the linear function  $\mathbf{l}$  as follows,

(**Spread Function:**) For  $1/2 > \alpha > 0$ , let

$$c(\alpha) = \max \left( \sup_{\theta} \Pr_{\mathcal{A}} [\mathbf{l}(\mathcal{A}) \in [\theta - \alpha, \theta + \alpha]], \sup_{\theta} \Pr_{\mathcal{B}} [\mathbf{l}(\mathcal{B}) \in [\theta - \alpha, \theta + \alpha]] \right)$$

then for all  $\tilde{\theta}$ ,

$$\left| \mathbf{E}_{\mathcal{A}} [\text{pos}(\mathbf{l}(\mathcal{A}) - \tilde{\theta})] - \mathbf{E}_{\mathcal{B}} [\text{pos}(\mathbf{l}(\mathcal{B}) - \tilde{\theta})] \right| \leq O\left(\frac{1}{\alpha^4}\right) \sum_{i \in [R]} \|\mathbf{l}^{\{i\}}\|_1^4 + 2c(\alpha). \quad (13)$$

*Proof.* Let us prove equation (12) first. Let  $\mathcal{X}_i = \{\mathcal{B}^{\{1\}}, \dots, \mathcal{B}^{\{i-1\}}, \mathcal{B}^{\{i\}}, \mathcal{A}^{\{i+1\}}, \dots, \mathcal{A}^{\{R\}}\}$ .

We know that

$$\begin{aligned} \mathbf{E}_{\mathcal{A}} [\Psi(\mathbf{l}(\mathcal{A}) - \theta)] - \mathbf{E}_{\mathcal{B}} [\Psi(\mathbf{l}(\mathcal{B}) - \theta)] &= \mathbf{E}_{\mathcal{X}_0} [\Psi(\mathbf{l}(\mathcal{X}_0) - \theta)] - \mathbf{E}_{\mathcal{X}_R} [\Psi(\mathbf{l}(\mathcal{X}_R) - \theta)] \\ &= \sum_{i=1}^R \mathbf{E}_{\mathcal{X}_{i-1}} [\Psi(\mathbf{l}(\mathcal{X}_{i-1}) - \theta)] - \mathbf{E}_{\mathcal{X}_i} [\Psi(\mathbf{l}(\mathcal{X}_i) - \theta)]. \end{aligned}$$

Therefore, it suffices to prove

$$\left| \mathbf{E}_{\mathcal{X}_{i-1}} [\Psi(\mathbf{l}(\mathcal{X}_{i-1}) - \theta)] - \mathbf{E}_{\mathcal{X}_i} [\Psi(\mathbf{l}(\mathcal{X}_i) - \theta)] \right| \leq K \|\mathbf{l}^{\{i\}}\|_1^4. \quad (14)$$

Let  $\mathcal{Y}_i = \{\mathcal{B}^{\{1\}}, \dots, \mathcal{B}^{\{i-1\}}, \mathcal{A}^{\{i+1\}}, \dots, \mathcal{A}^{\{R\}}\}$  and we have  $\mathcal{X}_i = \{\mathcal{Y}_i, \mathcal{B}^{\{i\}}\}$  and  $\mathcal{X}_{i-1} = \{\mathcal{Y}_i, \mathcal{A}^{\{i\}}\}$ . Then

$$\mathbf{E}_{\mathcal{X}_{i-1}} [\Psi(\mathbf{l}(\mathcal{X}_{i-1}) - \theta)] - \mathbf{E}_{\mathcal{X}_i} [\Psi(\mathbf{l}(\mathcal{X}_i) - \theta)] = \mathbf{E}_{\mathcal{Y}_i} \left[ \mathbf{E}_{\mathcal{A}^{\{i\}}} [\Psi(\mathbf{l}(\mathcal{X}_{i-1}) - \theta)] - \mathbf{E}_{\mathcal{B}^{\{i\}}} [\Psi(\mathbf{l}(\mathcal{X}_i) - \theta)] \right]. \quad (15)$$

Notice that

$$\mathbf{l}(\mathcal{X}_{i-1}) - \theta = \langle \mathbf{l}^{\{i\}}, \mathcal{A}^{\{i\}} \rangle + \sum_{1 \leq j \leq i-1} \langle \mathbf{l}^{\{j\}}, \mathcal{B}^{\{j\}} \rangle + \sum_{i+1 \leq j \leq R} \langle \mathbf{l}^{\{j\}}, \mathcal{A}^{\{j\}} \rangle - \theta$$

and

$$\mathbf{l}(\mathcal{X}_i) - \theta = \langle \mathbf{l}^{\{i\}}, \mathcal{B}^{\{i\}} \rangle + \sum_{1 \leq j \leq i-1} \langle \mathbf{l}^{\{j\}}, \mathcal{B}^{\{j\}} \rangle + \sum_{i+1 \leq j \leq R} \langle \mathbf{l}^{\{j\}}, \mathcal{A}^{\{j\}} \rangle - \theta.$$

Take  $\theta' = \sum_{1 \leq j \leq i-1} \langle \mathbf{l}^{\{j\}}, \mathbf{B}^{\{j\}} \rangle + \sum_{i+1 \leq j \leq R} \langle \mathbf{l}^{\{j\}}, \mathbf{A}^{\{j\}} \rangle - \theta$ , We can further rewrite equation (15) as

$$\mathbf{E}_{\mathcal{Y}_i} \left[ \mathbf{E}_{\mathbf{A}^{\{i\}}} [\Psi(\langle \mathbf{l}^{\{i\}}, \mathbf{A}^{\{i\}} \rangle + \theta')] - \mathbf{E}_{\mathbf{B}^{\{i\}}} [\Psi(\langle \mathbf{l}^{\{i\}}, \mathbf{B}^{\{i\}} \rangle + \theta')] \right]. \quad (16)$$

Using the Taylor expansion of  $\Psi$ , we have that the inner expectation of equation (16) is equal to

$$\begin{aligned} & \left| \mathbf{E}_{\mathbf{A}^{\{i\}}} [\Psi(\theta') + \Psi'(\theta') \langle \mathbf{l}^{\{i\}}, \mathbf{A}^{\{i\}} \rangle + \frac{\Psi''(\theta')}{2} (\langle \mathbf{l}^{\{i\}}, \mathbf{A}^{\{i\}} \rangle)^2 + \frac{\Psi'''(\theta')}{6} (\langle \mathbf{l}^{\{i\}}, \mathbf{A}^{\{i\}} \rangle)^3 + \frac{\Psi''''(\delta_1)}{24} (\langle \mathbf{l}^{\{i\}}, \mathbf{A}^{\{i\}} \rangle)^4] \right. \\ & \left. - \mathbf{E}_{\mathbf{B}^{\{i\}}} [\Psi(\theta') + \Psi'(\theta') \langle \mathbf{l}^{\{i\}}, \mathbf{B}^{\{i\}} \rangle + \frac{\Psi''(\theta')}{2} (\langle \mathbf{l}^{\{i\}}, \mathbf{B}^{\{i\}} \rangle)^2 + \frac{\Psi'''(\theta')}{6} (\langle \mathbf{l}^{\{i\}}, \mathbf{B}^{\{i\}} \rangle)^3 + \frac{\Psi''''(\delta_2)}{24} (\langle \mathbf{l}^{\{i\}}, \mathbf{B}^{\{i\}} \rangle)^4] \right|. \end{aligned} \quad (17)$$

for some  $\delta_1, \delta_2 \in \mathbb{R}$ .

Using the fact that  $\mathbf{A}^{\{i\}}$  and  $\mathbf{B}^{\{i\}}$  have matching moments up to degree 3, we can upper bound equation (17) by

$$\left| \mathbf{E}_{\mathbf{A}^{\{i\}}} \left[ \frac{\Psi''''(\delta_1)}{24} (\langle \mathbf{l}^{\{i\}}, \mathbf{A}^{\{i\}} \rangle)^4 \right] - \mathbf{E}_{\mathbf{B}^{\{i\}}} \left[ \frac{\Psi''''(\delta_2)}{24} (\langle \mathbf{l}^{\{i\}}, \mathbf{B}^{\{i\}} \rangle)^4 \right] \right| \leq \frac{K}{12} \|\mathbf{l}^{\{i\}}\|_1^4.$$

In the last inequality, we use the fact that  $\Psi$  is  $K$ -bounded and  $\langle \mathbf{l}^{\{i\}}, \mathbf{A}^{\{i\}} \rangle \leq \|\mathbf{l}^{\{i\}}\|_1$ ,  $\langle \mathbf{l}^{\{i\}}, \mathbf{B}^{\{i\}} \rangle \leq \|\mathbf{l}^{\{i\}}\|_1$  since all random variables in  $\mathcal{A}, \mathcal{B}$  are bounded by 1.

Overall, we bound the inner expectation of equation (16) by  $\frac{K}{12} \|\mathbf{l}^{\{i\}}\|_1^4$ . This implies equation (16) and therefore equation (14) is bounded by  $\frac{K}{12} \|\mathbf{l}^{\{i\}}\|_1^4$ , establishing equation (12).

To prove equation (13), we need to use the following lemma.

**Lemma C.1.** (*[37], Lemma 3.21*) *There exists an absolute constant  $C$  such that  $\forall 0 < \lambda < \frac{1}{2}$ , there exists  $\frac{C}{\lambda^4}$ -bounded function  $\Phi_\lambda : \mathbb{R} \rightarrow [0, 1]$  which approximates the  $\text{pos}(x)$  function in the following sense:  $\Phi_\lambda(t) = 1$  for all  $t > \lambda$ ;  $\Phi_\lambda(t) = 0$  for  $t < -\lambda$ .*

By the above lemma, we can find a  $\frac{C}{\alpha^4}$ -bounded function  $\Phi_\alpha$  such that  $\Phi_\alpha(\mathbf{l}(\mathcal{A}) - \theta)$  is equal to  $\text{pos}(\mathbf{l}(\mathcal{A}) - \theta)$  except when  $\mathbf{l}(\mathcal{A}) \in [\theta - \alpha, \theta + \alpha]$  and  $\Phi_\alpha(\mathbf{l}(\mathcal{B}) - \theta)$  is equal to  $\text{pos}(\mathbf{l}(\mathcal{B}) - \theta)$  except when  $\mathbf{l}(\mathcal{B}) \in [\theta - \alpha, \theta + \alpha]$ . Also for any  $x \in \mathbb{R}$ ,  $|\text{pos}(x) - \Phi_\alpha(x)| \leq 1$  as  $\text{pos}(x)$  and  $\Phi_\alpha(x)$  are both in  $[0, 1]$ .

Overall, we have

$$\begin{aligned} & \left| \mathbf{E}_{\mathcal{A}} [\text{pos}(\mathbf{l}(\mathcal{A}) - \theta)] - \mathbf{E}_{\mathcal{B}} [\text{pos}(\mathbf{l}(\mathcal{B}) - \theta)] \right| \leq \left| \mathbf{E}_{\mathcal{A}} [\text{pos}(\mathbf{l}(\mathcal{A}) - \theta)] - \mathbf{E}_{\mathcal{A}} [\Phi_\alpha(\mathbf{l}(\mathcal{A}) - \theta)] \right| \\ & \quad + \left| \mathbf{E}_{\mathcal{A}} [\Phi_\alpha(\mathbf{l}(\mathcal{A}) - \theta)] - \mathbf{E}_{\mathcal{B}} [\Phi_\alpha(\mathbf{l}(\mathcal{B}) - \theta)] \right| + \left| \mathbf{E}_{\mathcal{B}} [\Phi_\alpha(\mathbf{l}(\mathcal{B}) - \theta)] - \mathbf{E}_{\mathcal{B}} [\text{pos}(\mathbf{l}(\mathcal{B}) - \theta)] \right| \\ & \leq \frac{C}{\alpha^4} \sum_{i \in [R]} \|\mathbf{l}^{\{i\}}\|_1^4 + 2c(\alpha). \end{aligned}$$

□

## D Hardness of Smooth $k$ -Label Cover

First we state the bipartite smooth Label Cover given by Khot [30]. Our reduction is similar to the one in [19] but in addition requires proving the smoothness property.

**Definition D.1.** A Label Cover problem  $\mathcal{L}(G(W, V, E), M, N, \{\pi^{v,w} | (w, v) \in E\})$  consists of a bipartite graph  $G(V, W, E)$  with bipartition  $V$  and  $W$ .  $M, N$  are two positive integers such that  $M > N$ . There are projection functions  $\pi^{v,w} : [M] \rightarrow [N]$  associated with each edge  $(w, v) \in E$  where  $v \in V, w \in W$ . All vertices in  $W$  have the same degree (i.e.,  $W$ -side regular). For any labeling  $\Lambda : V \rightarrow [M]$  and  $\Lambda : W \rightarrow [N]$ , an edge is said to be satisfied if  $\pi^{v,w}(\Lambda(v)) = \Lambda(w)$ . We define  $\text{Opt}(\mathcal{L})$  to be the maximum fraction of edges satisfied by any labeling.

**Theorem D.2.** There is an absolute constant  $\gamma > 0$  such that for all integer parameters  $u$  and  $J$ , it is NP-hard to distinguish the following two cases: A Label Cover problem  $\mathcal{L}(G(W, V, E), N, M, \{\pi^{v,w} | (w, v) \in E\})$  with  $M = 7^{(J+1)u}$  and  $N = 2^u 7^{Ju}$  having

- $\text{Opt}(\mathcal{L}) = 1$  or
- $\text{Opt}(\mathcal{L}) \leq 2^{-2\gamma u}$ .

In addition, the Label Cover has the following properties:

- for each  $\pi^{v,w}$  and any  $i \in [N]$ , we have  $|(\pi^{v,w})^{-1}(i)| \leq 4^u$ ;
- for a fixed vertex  $w$  and a randomly picked neighbor  $v$  of  $w$ ,

$$\forall i, j \in [M], \Pr[\pi^{v,w}(i) = \pi^{v,w}(j)] \leq 1/J.$$

Below we prove Theorem 5.1.

*Proof.* Given an instance of bipartite Label Cover  $\mathcal{L}(G(V, W, E), M, N, \{\pi^{v,w} | (w, v) \in E\})$ , we can convert it to a smooth  $k$ -LABEL COVER instance  $\mathcal{L}'$  as follows. The vertex set of  $\mathcal{L}'$  is  $V$  and we generate the hyperedge set  $E'$  and projections associated with the hyperedges in the following way:

1. pick a vertex  $w \in W$ ;
2. pick a  $k$ -tuple of  $w$ 's neighbors  $v_1, \dots, v_k$  and add a hyperedge  $e = (v_1, \dots, v_k)$  to  $E'$  with projections  $\pi^{v_i, e} = \pi^{v_i, w}$  for each  $i \in [k]$ .

**Completeness:** If  $\text{Opt}(\mathcal{L}) = 1$ , then there exists a labeling  $\Lambda$  such that for every edge  $(w, v) \in E$ ,  $\pi^{v,w}(\Lambda(v)) = \Lambda(w)$ . We can simply take the restriction of labeling  $\Lambda$  on  $V$  for the smooth  $k$ -LABEL COVER instance  $\mathcal{L}'$ . For any hyperedge  $e = (v_1, v_2, \dots, v_k)$  generated by  $w \in W$ , we know  $\pi^{v_i, e}(\Lambda(v_i)) = \Lambda(w) = \pi^{v_j, e}(\Lambda(v_j))$  for any  $i, j \in [k]$ .

**Soundness:** If  $Opt(\mathcal{L}) \leq 2^{-2\gamma u}$ , then we can weakly satisfy at most  $2k^2 2^{-\gamma u}$ -fraction of the hyperedges in  $\mathcal{L}'$ . This can be proved via contrapositive argument. Suppose there is a labeling strategy  $\Lambda$  (defined on  $V$ ) for the smooth  $k$ -LABEL COVER that weakly satisfies  $\alpha \geq 2k^2 2^{-\gamma u}$  fraction of the hyperedges. Extend the labelling to  $W$  as follows: For each vertex  $w \in W$  and a neighbor  $v \in V$ , let  $\pi_{v,w}(\Lambda(v))$  be the label *recommended* by  $v$  to  $w$ . Simply assign for every vertex  $w \in W$ , the label most *recommended* by its neighbours.

By the fact that  $\Lambda$  weakly satisfies  $\alpha$ -fraction of hyperedges in  $\mathcal{L}'$ , we know that if we pick a vertex  $w$  and randomly pick two of its neighbors  $v_1, v_2$  then

$$\Pr[\pi^{v_1,w}(\Lambda(v_1)) = \pi^{v_2,w}(\Lambda(v_2))] \geq \frac{\alpha}{\binom{k}{2}} \geq \frac{2\alpha}{k^2}.$$

By an averaging argument, at least  $\frac{\alpha}{k^2}$ -fraction of the vertices  $w \in W$ , will have the following property: among all the possible pairs of  $w$ 's neighbors, at least  $\frac{\alpha}{k^2}$ -fraction of pairs *recommend* the same label for  $w$ . Let us call such a  $w$  to be a *nice*. It is easy to see that for every *nice*  $w$ , the most *recommended* label is actually recommended by at least  $\frac{\alpha}{k^2}$  fraction of its neighbours. Hence, the extended labelling satisfies at least  $\alpha/k^2$  fraction of edges incident at each *nice*  $w \in W$ . Using  $W$ -side regularity, we conclude that the extended labelling satisfies  $\frac{\alpha^2}{k^4} = 4 \cdot 2^{-2\gamma u}$ -fraction the edges of  $\mathcal{L}$  – a contradiction.

**Smoothness of  $\mathcal{L}'$ :** For any given vertex  $v$  in  $\mathcal{L}'$ , we want so show that if we randomly pick an hyperedge  $e'$  containing  $v$ , then for the projection  $\pi^{v,e'}$  as defined in  $\mathcal{L}'$ ,

$$\forall i, j \in [M], \Pr[\pi^{v,e'}(i) = \pi^{v,e'}(j)] \leq \frac{1}{J}.$$

To see this, notice that all vertices in  $W$  have the same degree; picking a projection  $\pi^{v,e'}$  using the above procedure is the same as randomly picking a neighbor  $w$  of  $v$  and using the projection  $\pi^{v,w}$  defined in  $\mathcal{L}$ . Therefore,

$$\forall i, j \in [M], \Pr[\pi^{v,e'}(i) = \pi^{v,e'}(j)] = \Pr[\pi^{v,w}(i) = \pi^{v,w}(j)] \leq \frac{1}{J}.$$

□