
List-decodable Linear Regression

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 We give the first polynomial-time algorithm for robust regression in the list-
2 decodable setting where an adversary can corrupt a greater than $1/2$ fraction
3 of examples.
4 For any $\alpha < 1$, our algorithm takes as input a sample $\{(x_i, y_i)\}_{i \leq n}$ of n linear
5 equations where αn of the equations satisfy $y_i = \langle x_i, \ell^* \rangle + \zeta$ for some small noise
6 ζ and $(1 - \alpha)n$ of the equations are *arbitrarily* chosen. It outputs a list L of size
7 $O(1/\alpha)$ - a fixed constant - that contains an ℓ that is close to ℓ^* .
8 Our algorithm succeeds whenever the inliers are chosen from a *certifiably* anti-
9 concentrated distribution D . In particular, this gives a $(d/\alpha)^{O(1/\alpha^8)}$ time algorithm
10 to find a $O(1/\alpha)$ size list when the inlier distribution is standard Gaussian. For
11 discrete product distributions that are anti-concentrated only in *regular* directions,
12 we give an algorithm that achieves similar guarantee under the promise that ℓ^* has
13 all coordinates of the same magnitude. To complement our result, we prove that the
14 anti-concentration assumption on the inliers is information-theoretically necessary.
15 To solve the problem we introduce a new framework for list-decodable learning
16 that strengthens the “identifiability to algorithms” paradigm based on the sum-of-
17 squares method.

18 1 Introduction

19 In this work, we design algorithms for the problem of linear regression that are robust to training sets
20 with an overwhelming ($\gg 1/2$) fraction of adversarially chosen outliers.

21 Outlier-robust learning algorithms have been extensively studied (under the name *robust statistics*)
22 in mathematical statistics [54, 45, 31, 29]. However, the algorithms resulting from this line of work
23 usually run in time exponential in the dimension of the data [7]. An influential line of recent work
24 [35, 1, 18, 39, 9, 36, 37, 30, 16, 19, 34] has focused on designing *efficient* algorithms for outlier-robust
25 learning.

26 Our work extends this line of research. Our algorithms work in the “list-decodable learning” frame-
27 work. In this model, a majority of the training data (a $1 - \alpha$ fraction) can be adversarially corrupted
28 leaving only an $\alpha \ll 1/2$ fraction of “inliers”. Since uniquely recovering the underlying parameters
29 is information-theoretically *impossible* in such a setting, the goal is to output a list (with an absolute
30 constant size) of parameters, one of which matches the ground truth. This model was introduced
31 in [3] to give a discriminative framework for clustering. More recently, beginning with [9], various
32 works [20, 36] have considered this as a model of “untrusted” data.

33 There has been phenomenal progress in developing techniques for outlier-robust learning with a
34 *small* ($\ll 1/2$)-fraction of outliers (e.g. outlier “filters” [15, 16, 11, 17], separation oracles for
35 inliers [15] or the *sum-of-squares* method [37, 30, 36, 34]). In contrast, progress on algorithms that
36 tolerate the significantly harsher conditions in the list-decodable setting has been slower. The only

37 prior works [9, 20, 36] in this direction designed list-decodable algorithms for mean estimation via
38 problem-specific methods.

39 In this paper, we develop a principled technique to give the first efficient list-decodable learning
40 algorithm for the fundamental problem of *linear regression*. Our algorithm takes a corrupted set
41 of linear equations with an $\alpha \ll 1/2$ fraction of inliers and outputs a $O(1/\alpha)$ -size list of linear
42 functions, one of which is guaranteed to be close to the ground truth (i.e., the linear function that
43 correctly labels the inliers). A key conceptual insight in this result is that list-decodable regression
44 information-theoretically requires the inlier-distribution to be “anti-concentrated”. Our algorithm
45 succeeds whenever the distribution satisfies a stronger “certifiable anti-concentration” condition that
46 is algorithmically “usable”. This class includes the standard gaussian distribution and more generally,
47 any spherically symmetric distribution with strictly sub-exponential tails.

48 Prior to our work¹, the state-of-the-art outlier-robust algorithms for linear regression [34, 21, 14,
49 49] could handle only a small (< 0.1)-fraction of outliers even under strong assumptions on the
50 underlying distributions.

51 List-decodable regression generalizes the well-studied [12, 32, 26, 57, 2, 10, 58, 51, 42] and *easier*
52 problem of *mixed linear regression*: given k “clusters” of examples that are labeled by one out of k
53 distinct unknown linear functions, find the unknown set of linear functions. All known techniques
54 for the problem rely on faithfully estimating certain *moment tensors* from samples and thus, cannot
55 tolerate the overwhelming fraction of outliers in the list-decodable setting. On the other hand, since
56 we can take any cluster as inliers and treat rest as outliers, our algorithm immediately yields new
57 efficient algorithms for mixed linear regression. Unlike all prior works, our algorithms work without
58 any pairwise separation or bounded condition-number assumptions on the k linear functions.

59 **List-Decodable Learning via the Sum-of-Squares Method** Our algorithm relies on a strengthen-
60 ing of the robust-estimation framework based on the sum-of-squares (SoS) method. This paradigm
61 has been recently used for clustering mixture models [30, 36] and obtaining algorithms for moment
62 estimation [37] and linear regression [34] that are resilient to a small ($\ll 1/2$) fraction of outliers
63 under the mildest known assumptions on the underlying distributions. At the heart of this technique is
64 a reduction of outlier-robust algorithm design to just finding “simple” proofs of unique “identifiability”
65 of the unknown parameter of the original distribution from a corrupted sample. However, this princi-
66 pled method works only in the setting with a small ($\ll 1/2$) fraction of outliers. As a consequence,
67 the work of [36] for mean estimation in the list-decodable setting relied on “supplementing” the SoS
68 method with a somewhat problem-dependent technique.

69 As an important conceptual contribution, our work yields a framework for list-decodable learning
70 that recovers some of the simplicity of the general blueprint. Central to our framework is a general
71 method of *rounding by votes* for “pseudo-distributions” in the setting with $\gg 1/2$ fraction outliers.
72 Our rounding builds on the work of [38] who developed such a method to give a simpler proof of the
73 list-decodable mean estimation result of [36]. In Section 2, we explain our ideas in detail.

74 The results in all the works above hold for any underlying distribution that has upper-bounded low-
75 degree moments and such bounds are “captured” within the SoS system. Such conditions are called as
76 “certified bounded moment” inequalities. An important contribution of this work is to formalize *anti-*
77 *concentration* inequalities within the SoS system and prove “certified anti-concentration” for natural
78 distribution families. Unlike bounded moment inequalities, there is no canonical encoding within
79 SoS for such statements. We choose an encoding that allow proving certified anti-concentration for a
80 distribution by showing the existence of a certain approximating polynomial. This allows showing
81 certified anti-concentration of natural distributions via a completely modular approach that relies on a
82 beautiful line of works that construct “weighted” polynomial approximators [43].

83 We believe that our framework for list-decodable estimation and our formulation of certified anti-
84 concentration condition will likely have further applications in outlier-robust learning.

85 1.1 Our Results

86 We first define our model for generating samples for list-decodable regression.

¹There’s a long line of work on robust regression algorithms (see for e.g. [8, 33]) that can tolerate corruptions only in the *labels*. We are interested in algorithms robust against corruptions in both examples and labels.

87 **Model 1.1** (Robust Linear Regression). For $0 < \alpha < 1$ and $\ell^* \in \mathbb{R}^d$ with $\|\ell^*\|_2 \leq 1$, let $\text{Lin}_D(\alpha, \ell^*)$
 88 denote the following probabilistic process to generate n noisy linear equations $\mathcal{S} = \{\langle x_i, a \rangle = y_i \mid$
 89 $1 \leq i \leq n\}$ in variable $a \in \mathbb{R}^d$ with αn inliers \mathcal{I} and $(1 - \alpha)n$ outliers \mathcal{O} :

- 90 1. Construct \mathcal{I} by choosing αn i.i.d. samples $x_i \sim D$ and set $y_i = \langle x_i, \ell^* \rangle + \zeta$ for additive
 91 noise ζ ,
- 92 2. Construct \mathcal{O} by choosing the remaining $(1 - \alpha)n$ equations arbitrarily and potentially
 93 adversarially w.r.t the inliers \mathcal{I} .

94 Note that α measures the “signal” (fraction of inliers) and can be $\ll 1/2$. The bound on the norm of
 95 ℓ^* is without any loss of generality. For the sake of exposition, we will restrict to $\zeta = 0$ for most of
 96 this paper and discuss (see Remarks 1.6 and 4.4) how our algorithms can tolerate additive noise.

97 An η -approximate algorithm for list-decodable regression takes input a sample from $\text{Lin}_D(\alpha, \ell^*)$ and
 98 outputs a *constant* (depending only on α) size list L of linear functions such that there is some $\ell \in L$
 99 that is η -close to ℓ^* .

100 One of our key conceptual contributions is to identify the strong relationship between *anti-*
 101 *concentration inequalities* and list-decodable regression. Anti-concentration inequalities are well-
 102 studied [22, 53, 50] in probability theory and combinatorics. The simplest of these inequalities upper
 103 bound the probability that a high-dimensional random variable has zero projections in any direction.

104 **Definition 1.2** (Anti-Concentration). A \mathbb{R}^d -valued zero-mean random variable Y has a δ -*anti-*
 105 *concentrated* distribution if $\Pr[\langle Y, v \rangle = 0] < \delta$.

106 In Proposition 2.4, we provide a simple but conceptually illuminating proof that anti-concentration is
 107 *sufficient* for list-decodable regression. In Theorem 6.1, we prove a sharp converse and show that
 108 anti-concentration is information-theoretically *necessary* for even noiseless list-decodable regression.
 109 This lower bound surprisingly holds for a natural distribution: uniform distribution on $\{0, 1\}^d$ and
 110 more generally, uniform distribution on $[q]^d$ for $q = \{0, 1, 2, \dots, q\}$. And in fact, our lower bound
 111 shows the impossibility of even the “easier” problem of mixed linear regression on this distribution.

112 **Theorem 1.3** (See Proposition 2.4 and Theorem 6.1). *There is a (inefficient) list-decodable regression*
 113 *algorithm for $\text{Lin}_D(\alpha, \ell^*)$ with list size $O(\frac{1}{\alpha})$ whenever D is α -anti-concentrated. Further, there*
 114 *exists a distribution D on \mathbb{R}^d that is $(\alpha + \epsilon)$ -anti-concentrated for every $\epsilon > 0$ but there is no*
 115 *algorithm for $\frac{\alpha}{2}$ -approximate list-decodable regression for $\text{Lin}_D(\alpha, \ell^*)$ that returns a list of size $< d$.*
 116

117 To handle additive noise of variance ζ^2 , we need a control of $\Pr[|\langle x, v \rangle| \leq \zeta]$. For our efficient
 118 algorithms, in addition, we need a *certified* version of the anti-concentration condition. Intuitively,
 119 certified anti-concentration asks for a *certificate* of the anti-concentration property of a random
 120 variable Y in the “sum-of-squares” proof system (see Section 3 for precise definitions). SoS is a
 121 proof system that reasons about polynomial inequalities. Since the “core indicator” $\mathbf{1}(|\langle x, v \rangle| \leq \delta)$
 122 is not a polynomial, we phrase the condition in terms of an approximating polynomial p . For this
 123 section, we will use “low-degree sum-of-squares proof” informally and encourage the reader to think
 124 of certified anti-concentration as a stronger version of anti-concentration that the SoS method can
 125 reason about.

126 **Definition 1.4** (Certifiable Anti-Concentration). A random variable Y has a k -*certifiably* (C, δ) -*anti-*
 127 *concentrated* distribution if there is a univariate polynomial p satisfying $p(0) = 1$ such that there is a
 128 degree k sum-of-squares proof of the following two inequalities:

- 129 1. $\forall v, \langle Y, v \rangle^2 \leq \delta^2 \mathbb{E} \langle Y, v \rangle^2$ implies $(p(\langle Y, v \rangle) - 1)^2 \leq \delta^2$.
- 130 2. $\forall v, \|v\|_2^2 \leq 1$ implies $\mathbb{E} p^2(\langle Y, v \rangle) \leq C\delta$.

131 We are now ready to state our main result.

132 **Theorem 1.5** (List-Decodable Regression). *For every $\alpha, \eta > 0$ and a k -certifiably $(C, \alpha^2 \eta^2 / 10C)$ -*
 133 *anti-concentrated distribution D on \mathbb{R}^d , there exists an algorithm that takes input a sample generated*
 134 *according to $\text{Lin}_D(\alpha, \ell^*)$ and outputs a list L of size $O(1/\alpha)$ such that there is an $\ell \in L$ satisfying*

Please note that sections 3-6 are in the supplementary material.

135 $\|\ell - \ell^*\|_2 < \eta$ with probability at least 0.99 over the draw of the sample. The algorithm needs a
 136 sample of size $n = (kd)^{O(k)}$ and runs in time $n^{O(k)} = (kd)^{O(k^2)}$.

137 **Remark 1.6** (Tolerating Additive Noise). For additive noise (not necessarily independent across
 138 samples) of variance ζ^2 in the inlier labels, our algorithm, in the same running time and sample
 139 complexity, outputs a list of size $O(1/\alpha)$ that contains an ℓ satisfying $\|\ell - \ell^*\|_2 \leq \frac{\zeta}{\alpha} + \eta$. Since we
 140 normalize ℓ^* to have unit norm, this guarantee is meaningful only when $\zeta \ll \alpha$.

141 **Remark 1.7** (Exponential Dependence on $1/\alpha$). List-decodable regression algorithms immediately
 142 yield algorithms for mixed linear regression (MLR) without any assumptions on the components.
 143 The state-of-the-art algorithms for MLR with gaussian components [42, 51] has an exponential
 144 dependence on $k = 1/\alpha$ in the running time in the absence of strong pairwise separation or small
 145 condition number of the components. Liang and Liu [42] (see Page 10 of their paper) use the
 146 relationship to learning mixtures of k gaussians (with an $\exp(k)$ lower bound [46]) to note that
 147 there may not exist any algorithms with polynomial dependence on $1/\alpha$ for MLR and thus, also for
 148 list-decodable regression.

149 **Certifiably anti-concentrated distributions** In Section 5, we show certifiable anti-concentration
 150 of some well-studied families of distributions. This includes the standard gaussian distribution and
 151 more generally any anti-concentrated spherically symmetric distribution with strictly sub-exponential
 152 tails. We also show that simple operations such as scaling, applying well-conditioned linear transfor-
 153 mations and sampling preserve certifiable anti-concentration. This yields:

154 **Corollary 1.8** (List-Decodable Regression for Gaussian Inliers). *For every $\alpha, \eta > 0$ there's*
 155 *an algorithm for list-decodable regression for the model $\text{Lin}_D(\alpha, \ell^*)$ with $D = \mathcal{N}(0, \Sigma)$ with*
 156 *$\lambda_{\max}(\Sigma)/\lambda_{\min}(\Sigma) = O(1)$ that needs $n = (d/\alpha\eta)^{O(\frac{1}{\alpha^4\eta^4})}$ samples and runs in time $n^{O(\frac{1}{\alpha^4\eta^4})} =$*
 157 *$(d/\alpha\eta)^{O(\frac{1}{\alpha^8\eta^8})}$.*

158 We note that certifiably anti-concentrated distributions are more restrictive compared to the families of
 159 distributions for which the most general robust estimation algorithms work [37, 36, 34]. To a certain
 160 extent, this is inherent. The families of distributions considered in these prior works do not satisfy
 161 anti-concentration in general. And as we discuss in more detail in Section 2, anti-concentration is
 162 information-theoretically *necessary* (see Theorem 1.3) for list-decodable regression. This surprisingly
 163 rules out families of distributions that might appear natural and “easy”, for example, the uniform
 164 distribution on $\{0, 1\}^n$.

165 We rescue this to an extent for the special case when ℓ^* in the model $\text{Lin}(\alpha, \ell^*)$ is a “Boolean
 166 vector”, i.e., has all coordinates of equal magnitude. Intuitively, this helps because while the the
 167 uniform distribution on $\{0, 1\}^n$ (and more generally, any discrete product distribution) is badly
 168 anti-concentrated in sparse directions, they are well anti-concentrated [22] in the directions that are
 169 far from any sparse vectors.

170 As before, for obtaining efficient algorithms, we need to work with a *certified* version (see Defini-
 171 tion 4.5) of such a restricted anti-concentration condition. As a specific Corollary (see Theorem 4.6
 172 for a more general statement), this allows us to show:

173 **Theorem 1.9** (List-Decodable Regression for Hypercube Inliers). *For every $\alpha, \eta > 0$ there's an*
 174 *η -approximate algorithm for list-decodable regression for the model $\text{Lin}_D(\alpha, \ell^*)$ with D is uniform*
 175 *on $\{0, 1\}^d$ that needs $n = (d/\alpha\eta)^{O(\frac{1}{\alpha^4\eta^4})}$ samples and runs in time $n^{O(\frac{1}{\alpha^4\eta^4})} = (d/\alpha\eta)^{O(\frac{1}{\alpha^8\eta^8})}$.*

176 In Section 4.1, we obtain similar results for general product distributions. It is an important open
 177 problem to prove certified anti-concentration for a broader family of distributions.

178 2 Overview of our Technique

179 In this section, we give a bird’s eye view of our approach and illustrate the important ideas in our
 180 algorithm for list-decodable regression. Thus, given a sample $\mathcal{S} = \{(x_i, y_i)\}_{i=1}^n$ from $\text{Lin}_D(\alpha, \ell^*)$,
 181 we must construct a constant-size list L of linear functions containing an ℓ close to ℓ^* .

Please note that sections 3-6 are in the supplementary material.

182 Our algorithm is based on the sum-of-squares method. We build on the “identifiability to algorithms”
 183 paradigm developed in several prior works [5, 4, 44, 37, 30, 36, 34] with some important conceptual
 184 differences.

185 **An inefficient algorithm** Let’s start by designing an inefficient algorithm for the problem. This
 186 may seem simple at the outset. But as we’ll see, solving this relaxed problem will rely on some
 187 important conceptual ideas that will serve as a starting point for our efficient algorithm.

188 Without computational constraints, it is natural to just return the list L of all linear functions ℓ that
 189 correctly labels all examples in some $S \subseteq \mathcal{S}$ of size αn . We call such an S , a large, *soluble* set. True
 190 inliers \mathcal{I} satisfy our search criteria so $\ell^* \in L$. However, it’s not hard to show (Proposition B.1) that
 191 one can choose outliers so that the list so generated has size $\exp(d)$ (far from a fixed constant!).

192 A potential fix is to search instead for a *coarse soluble partition* of \mathcal{S} , if it exists, into disjoint
 193 S_1, S_2, \dots, S_k and linear functions $\ell_1, \ell_2, \dots, \ell_k$ so that every $|S_i| \geq \alpha n$ and ℓ_i correctly computes
 194 the labels in S_i . In this setting, our list is small ($k \leq 1/\alpha$). But it is easy to construct samples \mathcal{S} for
 195 which this fails because there are coarse soluble partitions of \mathcal{S} where every ℓ_i is far from ℓ^* .

196 **Anti-Concentration** It turns out that any (even inefficient) algorithm for list-decodable regression
 197 provably (see Theorem 6.1) *requires* that the distribution of inliers² be sufficiently *anti-concentrated*:

198 **Definition 2.1** (Anti-Concentration). A \mathbb{R}^d -valued random variable Y with mean 0 is δ -anti-
 199 concentrated³ if for all non-zero v , $\Pr[\langle Y, v \rangle = 0] < \delta$. A set $T \subseteq \mathbb{R}^d$ is δ -anti-concentrated
 200 if the uniform distribution on T is δ -anti-concentrated.

201 As we discuss next, anti-concentration is also *sufficient* for list-decodable regression. Intuitively,
 202 this is because anti-concentration of the inliers prevents the existence of a soluble set that intersects
 203 significantly with \mathcal{I} and yet can be labeled correctly by $\ell \neq \ell^*$. This is simple to prove in the special
 204 case when \mathcal{S} admits a coarse soluble partition.

205 **Proposition 2.2.** *Suppose \mathcal{I} is α -anti-concentrated. Suppose there exists a partition*
 206 *$S_1, S_2, \dots, S_k \subseteq \mathcal{S}$ such that each $|S_i| \geq \alpha n$ and there exist $\ell_1, \ell_2, \dots, \ell_k$ such that $y_j = \langle \ell_i, x_j \rangle$*
 207 *for every $j \in S_i$. Then, there is an i such that $\ell_i = \ell^*$.*

208 *Proof.* Since $k \leq 1/\alpha$, there is a j such that $|\mathcal{I} \cap S_j| \geq \alpha |\mathcal{I}|$. Then, $\langle x_i, \ell_j \rangle = \langle x_i, \ell^* \rangle$ for every
 209 $i \in \mathcal{I} \cap S_j$. Thus, $\Pr_{i \sim \mathcal{I}}[\langle x_i, \ell_j - \ell^* \rangle = 0] \geq \alpha$. This contradicts anti-concentration of \mathcal{I} unless
 210 $\ell_j - \ell^* = 0$. □

211 The above proposition allows us to use *any* soluble partition as a *certificate* of correctness for the
 212 associated list L . Two aspects of this certificate were crucial in the above argument: 1) *largeness*:
 213 each S_i is of size αn - so the generated list is small, and, 2) *uniformity*: every sample is used in
 214 exactly one of the sets so \mathcal{I} must intersect one of the S_i s in at least α -fraction of the points.

215 **Identifiability via anti-concentration** For arbitrary \mathcal{S} , a coarse soluble partition might not exist.
 216 So we will generalize coarse soluble partitions to obtain certificates that exist for every sample \mathcal{S}
 217 and guarantee largeness and a relaxation of uniformity (formalized below). For this purpose, it is
 218 convenient to view such certificates as distributions μ on $\geq \alpha n$ size soluble subsets of \mathcal{S} so any
 219 collection $\mathcal{C} \subseteq 2^{\mathcal{S}}$ of αn size sets corresponds to the uniform distribution μ on \mathcal{C} .

220 To precisely define uniformity, let $W_i(\mu) = \mathbb{E}_{S \sim \mu}[\mathbf{1}(i \in S)]$ be the “frequency of i ”, that is,
 221 probability that the i th sample is chosen to be in a set drawn according to μ . Then, the uniform
 222 distribution μ on any coarse soluble k -partition satisfies $W_i = \frac{1}{k}$ for every i . That is, all samples
 223 $i \in \mathcal{S}$ are *uniformly* used in such a μ . To generalize this idea, we define $\sum_i W_i(\mu)^2$ as the *distance*
 224 *to uniformity* of μ . Up to a shift, this is simply the variance in the frequencies of the points in \mathcal{S}
 225 used in draws from μ . Our generalization of a coarse soluble partition of \mathcal{S} is any μ that minimizes
 226 $\sum_i W_i(\mu)^2$, the distance to uniformity, and is thus *maximally uniform* among all distributions
 227 supported on large soluble sets. Such a μ can be found by convex programming.

Please note that sections 3-6 are in the supplementary material.

²As in the standard robust estimation setting, the outliers are arbitrary and potentially adversarially chosen.

³Definition 1.4 differs slightly to handle list-decodable regression with additive noise in the inliers.

228 The following claim generalizes Proposition 2.2 to derive the same conclusion starting from any
 229 maximally uniform distribution supported on large soluble sets.

230 **Proposition 2.3.** *For a maximally uniform μ on αn size soluble subsets of \mathcal{S} ,*
 231 $\sum_{i \in \mathcal{I}} \mathbb{E}_{S \sim \mu} [\mathbf{1}(i \in S)] \geq \alpha |\mathcal{I}|$.

232 The proof proceeds by contradiction (see Lemma 4.3). We show that if $\sum_{i \in \mathcal{I}} W_i(\mu) \leq \alpha |\mathcal{I}|$, then we
 233 can strictly reduce the distance to uniformity by taking a mixture of μ with the distribution that places
 234 all its probability mass on \mathcal{I} . This allow us to obtain an (inefficient) algorithm for list-decodable
 235 regression establishing identifiability.

236 **Proposition 2.4** (Identifiability for List-Decodable Regression). *Let S be sample from $\text{Lin}(\alpha, \ell^*)$*
 237 *such that \mathcal{I} is δ -anti-concentrated for $\delta < \alpha$. Then, there's an (inefficient) algorithm that finds a list*
 238 *L of size $\frac{20}{\alpha - \delta}$ such that $\ell^* \in L$ with probability at least 0.99.*

239 *Proof.* Let μ be any maximally uniform distribution over αn size soluble subsets of \mathcal{S} . For $k = \frac{20}{\alpha - \delta}$,
 240 let S_1, S_2, \dots, S_k be independent samples from μ . Output the list L of k linear functions that
 241 correctly compute the labels in each S_i .

242 To see why $\ell^* \in L$, observe that $\mathbb{E}|S_j \cap \mathcal{I}| = \sum_{i \in \mathcal{I}} \mathbb{E} \mathbf{1}(i \in S_j) \geq \alpha |\mathcal{I}|$. By averaging, $\Pr[|S_j \cap \mathcal{I}| \geq$
 243 $\frac{\alpha + \delta}{2} |\mathcal{I}|] \geq \frac{\alpha - \delta}{2}$. Thus, there's a $j \leq k$ so that $|S_j \cap \mathcal{I}| \geq \frac{\alpha + \delta}{2} |\mathcal{I}|$ with probability at least
 244 $1 - (1 - \frac{\alpha - \delta}{2})^{\frac{20}{\alpha - \delta}} \geq 0.99$. We can now repeat the argument in the proof of Proposition 2.2 to
 245 conclude that any linear function that correctly labels S_j must equal ℓ^* . \square

246 **An efficient algorithm** Our identifiability proof suggests the following simple algorithm: 1) find
 247 any maximally uniform distribution μ on soluble subsets of size αn of \mathcal{S} , 2) take $O(1/\alpha)$ samples
 248 S_i from μ and 3) return the list of linear functions that correctly label the equations in S_i s. This is
 249 inefficient because searching over distributions is NP-hard in general.

250 To make this into an efficient algorithm, we start by observing that soluble subsets $S \subseteq \mathcal{S}$ of size αn
 251 can be described by the following set of quadratic equations where w stands for the indicator of S
 252 and ℓ , the linear function that correctly labels the examples in S .

$$A_{w, \ell}: \left\{ \begin{array}{l} \sum_{i=1}^n w_i = \alpha n \\ \forall i \in [n]. \quad w_i^2 = w_i \\ \forall i \in [n]. \quad w_i \cdot (y_i - \langle x_i, \ell \rangle) = 0 \\ \|\ell\|^2 \leq 1 \end{array} \right\} \quad (2.1)$$

253 Our efficient algorithm searches for a maximally uniform *pseudo-distribution* on w satisfying (2.1).
 254 Degree k pseudo-distributions (see Section 3 for precise definitions) are generalization of distributions
 255 that nevertheless “behave” just as distributions whenever we take (pseudo)-expectations (denoted
 256 by $\tilde{\mathbb{E}}$) of a class of degree k polynomials. And unlike distributions, degree k pseudo-distributions
 257 satisfying⁴ polynomial constraints (such as (2.1)) can be computed in time $n^{O(k)}$.

258 For the sake of intuition, it might be helpful to (falsely) think of pseudo-distributions $\tilde{\mu}$ as simply
 259 distributions where we only get access to moments of degree $\leq k$. Thus, we are allowed to compute
 260 expectations of all degree $\leq k$ polynomials with respect to $\tilde{\mu}$. Since $W_i(\tilde{\mu}) = \tilde{\mathbb{E}}_{\tilde{\mu}} w_i$ are just
 261 first moments of $\tilde{\mu}$, our notion of maximally uniform distributions extends naturally to pseudo-
 262 distributions. This allows us to prove an analog of Proposition 2.3 for pseudo-distributions and gives
 263 us an efficient replacement for Step 1.

264 **Proposition 2.5.** *For any maximally uniform $\tilde{\mu}$ of degree ≥ 2 , $\sum_{i \in \mathcal{I}} \tilde{\mathbb{E}}_{\tilde{\mu}}[w_i] \geq \alpha |\mathcal{I}| =$
 265 $\alpha \sum_{i \in [n]} \tilde{\mathbb{E}}_{\tilde{\mu}}[w_i]$.*

266 For Step 2, however, we hit a wall: it's not possible to obtain independent samples from $\tilde{\mu}$ given only
 267 low-degree moments.

Please note that sections 3-6 are in the supplementary material.

⁴See Fact 3.3 for a precise statement.

268 **Rounding by Votes** To circumvent this hurdle, our algorithm departs from rounding strategies for
 269 pseudo-distributions used in prior works and instead “rounds” *each* sample to a candidate linear
 270 function. While a priori, this method produces n different candidates instead of one, we will be able
 271 to extract a list of $O(\frac{1}{\alpha})$ size that contains the true vector from them. This step will crucially rely on
 272 anti-concentration properties of \mathcal{I} .

273 Consider the vector $v_i = \frac{\tilde{\mathbb{E}}_{\tilde{\mu}}[w_i \ell]}{\tilde{\mathbb{E}}_{\tilde{\mu}}[w_i]}$ whenever $\tilde{\mathbb{E}}_{\tilde{\mu}}[w_i] \neq 0$ (set v_i to zero, otherwise). This is simply the
 274 (scaled) average, according to $\tilde{\mu}$, of all the linear functions ℓ that are used to label the sets S of size
 275 αn in the support of $\tilde{\mu}$ whenever $i \in S$. Further, v_i depends only on the first two moments of $\tilde{\mu}$.

276 We think of v_i s as “votes” cast by the i th sample for the unknown linear function. Let us focus
 277 our attention on the votes v_i of $i \in \mathcal{I}$ - the inliers. We will show that according to the distribution
 278 proportional to $\tilde{\mathbb{E}}[w]$, the average ℓ_2 distance of v_i from ℓ^* is at max η :

$$\frac{1}{\sum_{i \in \mathcal{I}} \tilde{\mathbb{E}}[w_i]} \sum_{i \in \mathcal{I}} \tilde{\mathbb{E}}[w_i] \|v_i - \ell^*\|_2 < \eta. \quad (\star)$$

279 Before diving into (\star) , let’s see how it gives us our efficient list-decodable regression algorithm:

- 280 1. Find a pseudo-distribution $\tilde{\mu}$ satisfying (2.1) that minimizes distance to uniformity
 281 $\sum_i \tilde{\mathbb{E}}_{\tilde{\mu}}[w_i]^2$.
- 282 2. For $O(\frac{1}{\alpha})$ times, independently choose a random index $i \in [n]$ with probability proportional
 283 to $\tilde{\mathbb{E}}_{\tilde{\mu}}[w_i]$ and return the list of corresponding v_i s.

284 Step 1 above is a convex program - it minimizes a norm subject on the convex set of pseudo-
 285 distributions - and can be solved in polynomial time. Let’s analyze step 2 to see why the algorithm
 286 works. Using (\star) and Markov’s inequality, conditioned on $i \in \mathcal{I}$, $\|v_i - \ell^*\|_2 \leq 2\eta$ with probability
 287 $\geq 1/2$. By Proposition 2.5, $\frac{\sum_{i \in \mathcal{I}} \tilde{\mathbb{E}}[w_i]}{\sum_{i \in [n]} \tilde{\mathbb{E}}[w_i]} \geq \alpha$ so $i \in \mathcal{I}$ with probability at least α . Thus in each
 288 iteration of step 2, with probability at least $\alpha/2$, we choose an i such that v_i is 2η -close to ℓ^* .
 289 Repeating $O(1/\alpha)$ times gives us the 0.99 chance of success.

290 (\star) **via anti-concentration** As in the information-theoretic argument, (\star) relies on the anti-
 291 concentration of \mathcal{I} . Let’s do a quick proof for the case when $\tilde{\mu}$ is an actual distribution μ .

292 *Proof of (\star) for actual distributions μ .* Observe that μ is a distribution over (w, ℓ) satisfying (2.1).
 293 Recall that w indicates a subset $S \subseteq \mathcal{S}$ of size αn and $w_i = 1$ iff $i \in S$. And $\ell \in \mathbb{R}^d$ satisfies all the
 294 equations in S .

295 By Cauchy-Schwarz, $\sum_i \|\mathbb{E}_\mu[w_i \ell] - \mathbb{E}_\mu[w_i] \ell^*\| \leq \mathbb{E}_\mu[\sum_{i \in \mathcal{I}} w_i \|\ell - \ell^*\|]$. Next, as in Proposition 2.2,
 296 since \mathcal{I} is η -anti-concentrated, and for all S such that $|\mathcal{I} \cap S| \geq \eta |\mathcal{I}|$, $\ell - \ell^* = 0$. Thus, any such S
 297 in the support of μ contributes 0 to the expectation above. We will now show that the contribution
 298 from the remaining terms is upper bounded by η . Observe that since $\|\ell - \ell^*\| \leq 2$,
 299 $\mathbb{E}_\mu[\sum_{i \in \mathcal{I}} w_i \|\ell - \ell^*\|] = \mathbb{E}_\mu[\mathbf{1}(|S \cap \mathcal{I}| < \eta |\mathcal{I}|) \sum_{i \in S \cap \mathcal{I}} w_i \|\ell - \ell^*\|] \leq 2\eta |\mathcal{I}|$. \square

300 **SoSizing Anti-Concentration** The key to proving (\star) for pseudo-distributions is a *sum-of-squares*
 301 (SoS) proof of anti-concentration inequality: $\Pr_{x \sim \mathcal{I}}[\langle x, v \rangle = 0] \leq \eta$ in variable v . SoS is a restricted
 302 system for proving polynomial inequalities subject to polynomial inequality constraints. Thus, to
 303 even ask for a SoS proof we must phrase anti-concentration as a polynomial inequality.

304 To do this, let $p(z)$ be a low-degree polynomial approximator for the function $\mathbf{1}(z = 0)$.

305 Then, we can hope to “replace” the use of the inequality $\Pr_{x \sim \mathcal{I}}[\langle x, v \rangle = 0] \leq \eta \equiv \mathbb{E}_{x \sim \mathcal{I}}[\mathbf{1}(\langle x, v \rangle = 0)] \leq \eta$
 306 in the argument above by $\mathbb{E}_{x \sim \mathcal{I}}[p(\langle x, v \rangle)] \leq \eta$. Since polynomials grow unboundedly for
 307 large enough inputs, it is *necessary* for the uniform distribution on \mathcal{I} to have sufficiently light-tails
 308 to ensure that $\mathbb{E}_{x \sim \mathcal{I}} p(\langle x, v \rangle)$ is small. In Lemma A.1, we show that anti-concentration and strictly
 309 sub-exponential tails are *sufficient* to construct such a polynomial.

Please note that sections 3-6 are in the supplementary material.

310 We can finally ask for a SoS proof for $\mathbb{E}_{x \sim \mathcal{I}P}(\langle x, v \rangle) \leq \eta$ in variable v . We prove such *certified*
 311 anti-concentration inequalities for broad families of inlier distributions in Section 5.

312 3 Preliminaries

313 In this section, we define pseudo-distributions and sum-of-squares proofs. See the lecture notes [6]
 314 for more details and the appendix in [44] for proofs of the propositions appearing here.

315 Let $x = (x_1, x_2, \dots, x_n)$ be a tuple of n indeterminates and let $\mathbb{R}[x]$ be the set of polynomials
 316 with real coefficients and indeterminates x_1, \dots, x_n . We say that a polynomial $p \in \mathbb{R}[x]$ is a
 317 *sum-of-squares (sos)* if there are polynomials q_1, \dots, q_r such that $p = q_1^2 + \dots + q_r^2$.

318 3.1 Pseudo-distributions

319 Pseudo-distributions are generalizations of probability distributions. We can represent a discrete (i.e.,
 320 finitely supported) probability distribution over \mathbb{R}^n by its probability mass function $D: \mathbb{R}^n \rightarrow \mathbb{R}$
 321 such that $D \geq 0$ and $\sum_{x \in \text{supp}(D)} D(x) = 1$. Similarly, we can describe a pseudo-distribution by its
 322 mass function. Here, we relax the constraint $D \geq 0$ and only require that D passes certain low-degree
 323 non-negativity tests.

324 Concretely, a *level- ℓ pseudo-distribution* is a finitely-supported function $D: \mathbb{R}^n \rightarrow \mathbb{R}$ such that
 325 $\sum_x D(x) = 1$ and $\sum_x D(x) f(x)^2 \geq 0$ for every polynomial f of degree at most $\ell/2$. (Here, the
 326 summations are over the support of D .) A straightforward polynomial-interpolation argument shows
 327 that every level- ∞ -pseudo distribution satisfies $D \geq 0$ and is thus an actual probability distribution.
 328 We define the *pseudo-expectation* of a function f on \mathbb{R}^d with respect to a pseudo-distribution D ,
 329 denoted $\tilde{\mathbb{E}}_{D(x)} f(x)$, as

$$\tilde{\mathbb{E}}_{D(x)} f(x) = \sum_x D(x) f(x) . \quad (3.1)$$

330 The degree- ℓ moment tensor of a pseudo-distribution D is the tensor $\mathbb{E}_{D(x)}(1, x_1, x_2, \dots, x_n)^{\otimes \ell}$. In
 331 particular, the moment tensor has an entry corresponding to the pseudo-expectation of all monomials
 332 of degree at most ℓ in x . The set of all degree- ℓ moment tensors of probability distribution is a
 333 convex set. Similarly, the set of all degree- ℓ moment tensors of degree d pseudo-distributions is also
 334 convex. Key to the algorithmic utility of pseudo-distributions is the fact that while there can be no
 335 efficient separation oracle for the convex set of all degree- ℓ moment tensors of an actual probability
 336 distribution, there's a separation oracle running in time $n^{O(\ell)}$ for the convex set of the degree- ℓ
 337 moment tensors of all level- ℓ pseudodistributions.

338 **Fact 3.1** ([52, 48, 47, 40]). *For any $n, \ell \in \mathbb{N}$, the following set has a $n^{O(\ell)}$ -time weak separation*
 339 *oracle (in the sense of [28]):*

$$\left\{ \tilde{\mathbb{E}}_{D(x)}(1, x_1, x_2, \dots, x_n)^{\otimes d} \mid \text{degree-}d \text{ pseudo-distribution } D \text{ over } \mathbb{R}^n \right\} . \quad (3.2)$$

340 This fact, together with the equivalence of weak separation and optimization [28] allows us to
 341 efficiently optimize over pseudo-distributions (approximately)—this algorithm is referred to as the
 342 sum-of-squares algorithm.

343 The *level- ℓ sum-of-squares algorithm* optimizes over the space of all level- ℓ pseudo-distributions that
 344 satisfy a given set of polynomial constraints—we formally define this next.

345 **Definition 3.2** (Constrained pseudo-distributions). Let D be a level- ℓ pseudo-distribution over \mathbb{R}^n .
 346 Let $\mathcal{A} = \{f_1 \geq 0, f_2 \geq 0, \dots, f_m \geq 0\}$ be a system of m polynomial inequality constraints. We say
 347 that D *satisfies the system of constraints \mathcal{A} at degree r* , denoted $D \Big|_r \mathcal{A}$, if for every $S \subseteq [m]$ and
 348 every sum-of-squares polynomial h with $\deg h + \sum_{i \in S} \max\{\deg f_i, r\}$,

$$\tilde{\mathbb{E}}_D h \cdot \prod_{i \in S} f_i \geq 0 .$$

349 We write $D \models \mathcal{A}$ (without specifying the degree) if $D \models_0 \mathcal{A}$ holds. Furthermore, we say that $D \models_r \mathcal{A}$
350 holds *approximately* if the above inequalities are satisfied up to an error of $2^{-n^\ell} \cdot \|h\| \cdot \prod_{i \in S} \|f_i\|$,
351 where $\|\cdot\|$ denotes the Euclidean norm⁵ of the coefficients of a polynomial in the monomial basis.

352 We remark that if D is an actual (discrete) probability distribution, then we have $D \models \mathcal{A}$ if and only
353 if D is supported on solutions to the constraints \mathcal{A} .

354 We say that a system \mathcal{A} of polynomial constraints is *explicitly bounded* if it contains a constraint of
355 the form $\{\|x\|^2 \leq M\}$. The following fact is a consequence of [Fact 3.1](#) and [\[28\]](#),

356 **Fact 3.3** (Efficient Optimization over Pseudo-distributions). *There exists an $(n + m)^{O(\ell)}$ -time*
357 *algorithm that, given any explicitly bounded and satisfiable system⁶ \mathcal{A} of m polynomial constraints*
358 *in n variables, outputs a level- ℓ pseudo-distribution that satisfies \mathcal{A} approximately.*

359 3.2 Sum-of-squares proofs

360 Let f_1, f_2, \dots, f_r and g be multivariate polynomials in x . A *sum-of-squares proof* that the constraints
361 $\{f_1 \geq 0, \dots, f_m \geq 0\}$ imply the constraint $\{g \geq 0\}$ consists of polynomials $(p_S)_{S \subseteq [m]}$ such that

$$g = \sum_{S \subseteq [m]} p_S \cdot \prod_{i \in S} f_i. \quad (3.3)$$

362 We say that this proof has *degree* ℓ if for every set $S \subseteq [m]$, the polynomial $p_S \prod_{i \in S} f_i$ has degree at
363 most ℓ . If there is a degree ℓ SoS proof that $\{f_i \geq 0 \mid i \leq r\}$ implies $\{g \geq 0\}$, we write:

$$\{f_i \geq 0 \mid i \leq r\} \Big|_{\ell} \{g \geq 0\}. \quad (3.4)$$

364 Sum-of-squares proofs satisfy the following inference rules. For all polynomials $f, g: \mathbb{R}^n \rightarrow \mathbb{R}$ and
365 for all functions $F: \mathbb{R}^n \rightarrow \mathbb{R}^m$, $G: \mathbb{R}^n \rightarrow \mathbb{R}^k$, $H: \mathbb{R}^p \rightarrow \mathbb{R}^n$ such that each of the coordinates of
366 the outputs are polynomials of the inputs, we have:

$$\frac{\mathcal{A} \Big|_{\ell} \{f \geq 0, g \geq 0\}}{\mathcal{A} \Big|_{\ell} \{f + g \geq 0\}}, \frac{\mathcal{A} \Big|_{\ell} \{f \geq 0\}, \mathcal{A} \Big|_{\ell'} \{g \geq 0\}}{\mathcal{A} \Big|_{\ell + \ell'} \{f \cdot g \geq 0\}} \quad (\text{addition and multiplication})$$

$$\frac{\mathcal{A} \Big|_{\ell} B, B \Big|_{\ell'} C}{\mathcal{A} \Big|_{\ell \cdot \ell'} C} \quad (\text{transitivity})$$

$$\frac{\{F \geq 0\} \Big|_{\ell} \{G \geq 0\}}{\{F(H) \geq 0\} \Big|_{\ell \cdot \deg(H)} \{G(H) \geq 0\}}. \quad (\text{substitution})$$

367 Low-degree sum-of-squares proofs are sound and complete if we take low-level pseudo-distributions
368 as models.

369 Concretely, sum-of-squares proofs allow us to deduce properties of pseudo-distributions that satisfy
370 some constraints.

371 **Fact 3.4** (Soundness). *If $D \models_r \mathcal{A}$ for a level- ℓ pseudo-distribution D and there exists a sum-of-squares*
372 *proof $\mathcal{A} \Big|_{r'} B$, then $D \models_{r \cdot r' + r'} B$.*

373 If the pseudo-distribution D satisfies \mathcal{A} only approximately, soundness continues to hold if we require
374 an upper bound on the bit-complexity of the sum-of-squares $\mathcal{A} \Big|_{r'} B$ (number of bits required to
375 write down the proof).

376 In our applications, the bit complexity of all sum of squares proofs will be $n^{O(\ell)}$ (assuming that
377 all numbers in the input have bit complexity $n^{O(1)}$). This bound suffices in order to argue about
378 pseudo-distributions that satisfy polynomial constraints approximately.

⁵The choice of norm is not important here because the factor 2^{-n^ℓ} swamps the effects of choosing another norm.

⁶Here, we assume that the bitcomplexity of the constraints in \mathcal{A} is $(n + m)^{O(1)}$.

379 The following fact shows that every property of low-level pseudo-distributions can be derived by
 380 low-degree sum-of-squares proofs.

381 **Fact 3.5** (Completeness). *Suppose $d \geq r' \geq r$ and \mathcal{A} is a collection of polynomial constraints with
 382 degree at most r , and $\mathcal{A} \vdash \{\sum_{i=1}^n x_i^2 \leq B\}$ for some finite B .*

383 *Let $\{g \geq 0\}$ be a polynomial constraint. If every degree- d pseudo-distribution that satisfies $D \Big|_{r'} \mathcal{A}$
 384 also satisfies $D \Big|_{r'} \{g \geq 0\}$, then for every $\epsilon > 0$, there is a sum-of-squares proof $\mathcal{A} \Big|_d \{g \geq -\epsilon\}$.*

385 We will use the following Cauchy-Schwarz inequality for pseudo-distributions:

386 **Fact 3.6** (Cauchy-Schwarz for Pseudo-distributions). *Let f, g be polynomials of degree at most d in
 387 indeterminate $x \in \mathbb{R}^d$. Then, for any degree d pseudo-distribution $\tilde{\mu}$, $\tilde{\mathbb{E}}_{\tilde{\mu}}[fg] \leq \sqrt{\tilde{\mathbb{E}}_{\tilde{\mu}}[f^2]} \sqrt{\tilde{\mathbb{E}}_{\tilde{\mu}}[g^2]}$.*

388 The following fact is a simple corollary of the fundamental theorem of algebra:

389 **Fact 3.7**. *For any univariate degree d polynomial $p(x) \geq 0$ for all $x \in \mathbb{R}$, $\Big|_d \{p(x) \geq 0\}$.*

390 This can be extended to univariate polynomial inequalities over intervals of \mathbb{R} .

391 **Fact 3.8** (Fekete and Markov-Lukács, see [41]). *For any univariate degree d polynomial $p(x) \geq 0$
 392 for $x \in [a, b]$, $\{x \geq a, x \leq b\} \Big|_d \{p(x) \geq 0\}$.*

393 4 Algorithm for List-Decodable Robust Regression

394 In this section, we describe and analyze our algorithm for list-decodable regression and prove our
 395 first main result restated here.

396 **Theorem 1.5** (List-Decodable Regression). *For every $\alpha, \eta > 0$ and a k -certifiably $(C, \alpha^2 \eta^2 / 10C)$ -
 397 anti-concentrated distribution D on \mathbb{R}^d , there exists an algorithm that takes input a sample generated
 398 according to $\text{Lin}_D(\alpha, \ell^*)$ and outputs a list L of size $O(1/\alpha)$ such that there is an $\ell \in L$ satisfying
 399 $\|\ell - \ell^*\|_2 < \eta$ with probability at least 0.99 over the draw of the sample. The algorithm needs a
 400 sample of size $n = (kd)^{O(k)}$ and runs in time $n^{O(k)} = (kd)^{O(k^2)}$.*

401 We will analyze Algorithm 1 to prove Theorem 1.5.

$$402 \mathcal{A}_{w, \ell} : \left\{ \begin{array}{l} \sum_{i=1}^n w_i = \alpha n \\ \forall i \in [n]. \quad w_i^2 = w_i \\ \forall i \in [n]. \quad w_i \cdot (y_i - \langle x_i, \ell \rangle) = 0 \\ \sum_{i \leq d} \ell_i^2 \leq 1 \end{array} \right\} \quad (4.1)$$

Algorithm 1 (List-Decodable Regression).

Given: Sample \mathcal{S} of size n drawn according to $\text{Lin}(\alpha, n, \ell^*)$ with inliers \mathcal{I} , $\eta > 0$.

Output: A list $L \subseteq \mathbb{R}^d$ of size $O(1/\alpha)$ such that there exists a $\ell \in L$ satisfying $\|\ell - \ell^*\|_2 < \eta$.

Operation:

1. Find a degree $O(1/\alpha^4 \eta^4)$ pseudo-distribution $\tilde{\mu}$ satisfying $\mathcal{A}_{w, \ell}$ that minimizes $\|\tilde{\mathbb{E}}[w]\|_2$.
2. For each $i \in [n]$ such that $\tilde{\mathbb{E}}_{\tilde{\mu}}[w_i] > 0$, let $v_i = \frac{\tilde{\mathbb{E}}_{\tilde{\mu}}[w_i \ell]}{\tilde{\mathbb{E}}_{\tilde{\mu}}[w_i]}$. Otherwise, set $v_i = 0$.
3. Take J be a random multiset formed by union of $O(1/\alpha)$ independent draws of $i \in [n]$ with probability $\frac{\tilde{\mathbb{E}}[w_i]}{\alpha n}$.
4. Output $L = \{v_i \mid i \in J\}$ where $J \subseteq [n]$.

403 Our analysis follows the discussion in the overview. We start by formally proving (\star).

404 **Lemma 4.1.** For any $t \geq k$ and any \mathcal{S} so that $\mathcal{I} \subseteq \mathcal{S}$ is k -certifiably $(C, \alpha^2 \eta^2 / 4C)$ -anti-
 405 concentrated,

$$\mathcal{A}_{w,\ell} \Big|_{\frac{w,\ell}{t}} \left\{ \frac{1}{|\mathcal{I}|} \sum_{i \in \mathcal{I}} w_i \|\ell - \ell^*\|_2^2 \leq \frac{\alpha^2 \eta^2}{4} \right\}$$

406

407 *Proof.* We start by observing:

$$\mathcal{A}_{w,\ell} \Big|_{\frac{\ell}{2}} \|\ell - \ell^*\|_2^2 \leq 2.$$

408 Since \mathcal{I} is $(C, \alpha \eta / 2C)$ -anti-concentrated, there exists a univariate polynomial p such that $\forall i$:

$$\{w_i \langle x, \ell - \ell^* \rangle = 0\} \Big|_{\frac{k}{\ell}} \{p(\langle x_i, \ell - \ell^* \rangle) = 1\} \quad (4.2)$$

409 and

$$\{\|\ell\|^2 \leq 1\} \Big|_{\frac{k}{\ell}} \left\{ \frac{1}{|\mathcal{I}|} \sum_{i \in \mathcal{I}} p(\langle x_i, \ell - \ell^* \rangle)^2 \leq \frac{\alpha^2 \eta^2}{4} \right\} \quad (4.3)$$

410 Using (4.2), we have:

$$\mathcal{A}_{w,\ell} \Big|_{\frac{w,\ell}{t+2}} \{1 - p^2(\langle x_i, \ell - \ell^* \rangle) = 0\} \Big|_{\frac{w,\ell}{t+2}} \{1 - w_i p^2(\langle x_i, \ell - \ell^* \rangle) = 0\}$$

411 Using (4.3) and $\mathcal{A}_{w,\ell} \Big|_{\frac{w}{t}} \{w_i^2 = w_i\}$, we thus have:

$$\begin{aligned} \mathcal{A}_{w,\ell} \Big|_{\frac{w,\ell}{t+2}} \left\{ \frac{1}{|\mathcal{I}|} \sum_{i \in \mathcal{I}} w_i \|\ell - \ell^*\|_2^2 \right\} &= \frac{1}{|\mathcal{I}|} \sum_{i \in \mathcal{I}} w_i \|\ell - \ell^*\|_2^2 w_i p^2(\langle x_i, \ell - \ell^* \rangle) = \frac{1}{|\mathcal{I}|} \sum_{i \in \mathcal{I}} w_i \|\ell - \ell^*\|_2^2 p^2(\langle x_i, \ell - \ell^* \rangle) \\ &\leq \frac{1}{|\mathcal{I}|} \sum_{i \in \mathcal{I}} \|\ell - \ell^*\|_2^2 p^2(\langle x_i, \ell - \ell^* \rangle) \leq \frac{\alpha^2 \eta^2}{4}. \end{aligned}$$

412

□

413 As a consequence of this lemma, we can show that a constant fraction of the v_i for $i \in \mathcal{I}$ constructed
 414 in the algorithm are close to ℓ^* .

415 **Lemma 4.2.** For any $\tilde{\mu}$ of degree k satisfying $\mathcal{A}_{w,\ell}, \frac{1}{|\mathcal{I}|} \sum_{i \in \mathcal{I}} \tilde{\mathbb{E}}[w_i] \cdot \|v_i - \ell^*\|_2 \leq \frac{\alpha}{2} \eta$.

416 *Proof.* By Lemma 4.1, we have: $\mathcal{A}_{w,\ell} \Big|_{\frac{w,\ell}{k}} \left\{ \frac{1}{|\mathcal{I}|} \sum_{i \in \mathcal{I}} w_i \|\ell - \ell^*\|_2^2 \leq \frac{\alpha^2 \eta^2}{4} \right\}$.

417 We also have: $\mathcal{A}_{w,\ell} \Big|_{\frac{w,\ell}{2}} \{w_i^2 - w_i = 0\}$ for any i . This yields:

$$\mathcal{A}_{w,\ell} \Big|_{\frac{w,\ell}{k}} \left\{ \frac{1}{|\mathcal{I}|} \sum_{i \in \mathcal{I}} \|w_i \ell - w_i \ell^*\|_2^2 \leq \frac{\alpha^2 \eta^2}{4} \right\}$$

418 Since $\tilde{\mu}$ satisfies $\mathcal{A}_{w,\ell}$, taking pseudo-expectations yields: $\frac{1}{|\mathcal{I}|} \sum_{i \in \mathcal{I}} \tilde{\mathbb{E}} \|w_i \ell - w_i \ell^*\|_2^2 \leq \frac{\alpha^2 \eta^2}{4}$.

419 By Cauchy-Schwarz for pseudo-distributions (Fact 3.6), we have:

$$\left(\frac{1}{|\mathcal{I}|} \sum_{i \in \mathcal{I}} \|\tilde{\mathbb{E}}[w_i \ell] - \tilde{\mathbb{E}}[w_i] \ell^*\|_2 \right)^2 \leq \frac{1}{|\mathcal{I}|} \sum_{i \in \mathcal{I}} \|\tilde{\mathbb{E}}[w_i \ell] - \tilde{\mathbb{E}}[w_i] \ell^*\|_2^2 \leq \frac{\alpha^2 \eta^2}{4}.$$

420 Using $v_i = \frac{\tilde{\mathbb{E}}[w_i \ell]}{\tilde{\mathbb{E}}[w_i]}$ if $\tilde{\mathbb{E}}[w_i] > 0$ and 0 otherwise, we have: $\frac{1}{|\mathcal{I}|} \sum_{i \in \mathcal{I}, \tilde{\mathbb{E}}[w_i] > 0} \tilde{\mathbb{E}}[w_i] \cdot \|v_i - \ell^*\|_2 \leq \frac{\alpha}{2} \eta$.

421

□

422 Next, we formally prove that maximally uniform pseudo-distributions satisfy Proposition 2.5.

423 **Lemma 4.3.** For any $\tilde{\mu}$ of degree ≥ 4 satisfying $\mathcal{A}_{w,\ell}$ that minimizes $\|\tilde{\mathbb{E}}[w]\|_2$, $\sum_{i \in \mathcal{I}} \tilde{\mathbb{E}}_{\tilde{\mu}}[w_i] \geq \alpha^2 n$.
 424

425 *Proof.* Let $u = \frac{1}{\alpha n} \tilde{\mathbb{E}}[w]$. Then, u is a non-negative vector satisfying $\sum_{i \sim [n]} u_i = 1$.

426 Let $\text{wt}(\mathcal{I}) = \sum_{i \in \mathcal{I}} u_i$ and $\text{wt}(\mathcal{O}) = \sum_{i \notin \mathcal{I}} u_i$. Then, $\text{wt}(\mathcal{I}) + \text{wt}(\mathcal{O}) = 1$.

427 We will show that if $\text{wt}(\mathcal{I}) < \alpha$, then there's a pseudo-distribution $\tilde{\mu}'$ that satisfies $\mathcal{A}_{w,\ell}$ and has a
 428 lower value of $\|\tilde{\mathbb{E}}[w]\|_2$. This is enough to complete the proof.

429 To show this, we will “mix” $\tilde{\mu}$ with another pseudo-distribution satisfying $\mathcal{A}_{w,\ell}$. Let $\tilde{\mu}^*$ be the *actual*
 430 distribution supported on single (w, ℓ) - the indicator $\mathbf{1}_{\mathcal{I}}$ and ℓ^* . Thus, $\tilde{\mathbb{E}}_{\tilde{\mu}^*} w_i = 1$ iff $i \in \mathcal{I}$ and
 431 0 otherwise. $\tilde{\mu}^*$ clearly satisfies $\mathcal{A}_{w,\ell}$. Thus, any convex combination (mixture) of $\tilde{\mu}$ and $\tilde{\mu}^*$ also
 432 satisfies $\mathcal{A}_{w,\ell}$.

433 Let $\tilde{\mu}_\lambda = (1 - \lambda)\tilde{\mu} + \lambda\tilde{\mu}^*$. We will show that there is a $\lambda > 0$ such that $\|\tilde{\mathbb{E}}_{\tilde{\mu}_\lambda}[w]\|_2 < \|\tilde{\mathbb{E}}[w]\|_2$.

434 We first lower bound $\|u\|_2^2$ in terms of $\text{wt}(\mathcal{I})$ and $\text{wt}(\mathcal{O})$. Observe that for any fixed values of $\text{wt}(\mathcal{I})$
 435 and $\text{wt}(\mathcal{O})$, the minimum is attained by the vector u that ensures $u_i = \frac{1}{\alpha n} \text{wt}(\mathcal{I})$ for each $i \in \mathcal{I}$ and
 436 $u_i = \frac{1}{(1-\alpha)n} \text{wt}(\mathcal{O})$.

$$\text{This gives } \|u\|^2 \geq \left(\frac{\text{wt}(\mathcal{I})}{\alpha n}\right)^2 \alpha n + \left(\frac{1 - \text{wt}(\mathcal{I})}{(1 - \alpha)n}\right)^2 (1 - \alpha)n = \frac{1}{\alpha n} \cdot \left(\text{wt}(\mathcal{I}) + (1 - \text{wt}(\mathcal{I}))^2 \left(\frac{\alpha}{1 - \alpha}\right)\right).$$

437 Next, we compute the the ℓ_2 norm of $u' = \frac{1}{\alpha n} \tilde{\mathbb{E}}_{\tilde{\mu}_\lambda} w$ as:

$$\|u'\|_2^2 = (1 - \lambda)^2 \|u\|^2 + \frac{\lambda^2}{\alpha n} + 2\lambda(1 - \lambda) \frac{\text{wt}(\mathcal{I})}{\alpha n}.$$

438

$$\begin{aligned} \text{Thus, } \|u'\|^2 - \|u\|^2 &= (-2\lambda + \lambda^2) \|u\|^2 + \frac{\lambda^2}{\alpha n} + 2\lambda(1 - \lambda) \frac{\text{wt}(\mathcal{I})}{\alpha n} \\ &\leq \frac{-2\lambda + \lambda^2}{\alpha n} \cdot \left(\text{wt}(\mathcal{I})^2 + (1 - \text{wt}(\mathcal{I}))^2 \frac{\alpha}{1 - \alpha}\right) + \frac{\lambda^2}{\alpha n} + 2\lambda(1 - \lambda) \frac{\text{wt}(\mathcal{I})}{\alpha n} \end{aligned}$$

439

$$\begin{aligned} \text{Rearranging, } \|u\|^2 - \|u'\|^2 &\geq \frac{\lambda}{\alpha n} \left((2 - \lambda) \cdot \left(\text{wt}(\mathcal{I})^2 + (1 - \text{wt}(\mathcal{I}))^2 \left(\frac{\alpha}{1 - \alpha}\right)\right) - \lambda - 2(1 - \lambda)\text{wt}(\mathcal{I}) \right) \\ &\geq \frac{\lambda(2 - \lambda)}{\alpha n} \left(\text{wt}(\mathcal{I})^2 + (1 - \text{wt}(\mathcal{I}))^2 \frac{\alpha}{1 - \alpha} - \text{wt}(\mathcal{I}) \right) \end{aligned}$$

440 Now, whenever $\text{wt}(\mathcal{I}) < \alpha$, $\text{wt}(\mathcal{I})^2 + (1 - \text{wt}(\mathcal{I}))^2 \frac{\alpha}{1 - \alpha} - \text{wt}(\mathcal{I}) > 0$. Thus, we can choose a small
 441 enough $\lambda > 0$ so that $\|u\|^2 - \|u'\|^2 > 0$.

442

□

443 Lemma 4.3 and Lemma 4.2 immediately imply the correctness of our algorithm.

444 *Proof of Main Theorem 1.5.* First, since D is k -certifiably $(C, \alpha\eta/4C)$ -anti-concentrated,
 445 Lemma 5.5 implies taking $\geq n = (kd)^{O(k)}$ samples ensures that \mathcal{I} is k -certifiably $(C, \alpha\eta/2C)$ -anti-
 446 concentrated with probability at least $1 - 1/d$. Let's condition on this event in the following.

447 Let $\tilde{\mu}$ be a pseudo-distribution of degree t satisfying $\mathcal{A}_{w,\ell}$ and minimizing $\|\tilde{\mathbb{E}}[w]\|_2$. Such a pseudo-
 448 distribution exists as can be seen by just taking the distribution with a single-point support w where
 449 $w_i = 1$ iff $i \in \mathcal{I}$.

450 From Lemma 4.2, we have: $\frac{1}{|\mathcal{I}|} \sum_{i \in \mathcal{I}} \tilde{\mathbb{E}}[w_i] \cdot \|v_i - \ell^*\|_2 \leq \frac{\alpha}{2} \eta$. Let $Z = \frac{1}{\alpha n} \sum_{i \in \mathcal{I}} \tilde{\mathbb{E}}[w_i]$. By a
 451 rescaling, we obtain:

$$\frac{1}{|\mathcal{I}|} \sum_{i \in \mathcal{I}} \frac{\tilde{\mathbb{E}}[w_i]}{Z} \cdot \|v_i - \ell^*\|_2 \leq \frac{1}{Z} \frac{\alpha}{2} \eta. \quad (4.4)$$

452 Using Lemma 4.3, $Z \geq \alpha$. Thus,

$$\frac{1}{|\mathcal{I}|} \sum_{i \in \mathcal{I}} \frac{\tilde{\mathbb{E}}[w_i]}{Z} \cdot \|v_i - \ell^*\|_2 \leq \eta/2. \quad (4.5)$$

453 Let $i \in [n]$ be chosen with probability $\frac{\tilde{\mathbb{E}}[w_i]}{\alpha n}$. Then, $i \in \mathcal{I}$ with probability $Z \geq \alpha$. By Markov's
 454 inequality applied to (4.5), with $\frac{1}{2}$ conditioned on $i \in \mathcal{I}$, $\|v_i - \ell^*\|_2 < \eta$. Thus, in total, with
 455 probability at least $\alpha/2$, $\|v_i - \ell^*\|_2 \leq \eta$. Thus, the with probability at least 0.99 over the draw of the
 456 random set J , the list constructed by the algorithm contains an ℓ such that $\|\ell - \ell^*\|_2 \leq \eta$.

457 Let us now account for the running time and sample complexity of the algorithm. The sample
 458 size for the algorithm is dictated by Lemma 5.5 and is $(kd)^{O(k)}$, which for our choice of p goes
 459 as $(kd)^{O(k)}$. A pseudo-distribution satisfying $\mathcal{A}_{w,\ell}$ and minimizing $\|\tilde{\mathbb{E}}[w]\|_2$ can be found in time
 460 $n^{O(k)} = (kd)^{O(k^2)}$. The rounding procedure runs in time at most $O(nd)$. \square

461 *Remark 4.4* (Tolerating Additive Noise). To tolerate independent additive noise, our algorithm and
 462 analysis change minimally. For an additive noise of variance $\zeta^2 \ll \alpha^2 \eta^2$ in the inliers, we modify
 463 $\mathcal{A}_{w,\ell}$ by replacing the constraint $\forall i, w_i \cdot (y_i - \langle x_i, \ell \rangle) = 0$ by $\forall i, \pm w_i \cdot (y_i - \langle x_i, \ell \rangle) \leq 4\zeta$. And
 464 $\sum_{i=1}^n w_i = \alpha n$ to $\sum_{i=1}^n w_i = (\alpha/2)n$.

465 This means that instead of searching for a subsample of size αn that has a exact solution ℓ , we search
 466 for a subsample of size $\alpha/2n$ where there's a solution ℓ with an additive error of at most 2ζ . With
 467 additive noise of variance ζ^2 , it is easy to check that there's a subset of $1/2$ fraction of inliers that
 468 satisfies this property. Thus, $\mathcal{A}_{w,\ell}$ is feasible.

469 Our analysis remains exactly the same except for one change in the proof of Lemma 4.1. We start
 470 from a distribution that is $(C, \alpha\eta\zeta/100C)$ -certifiably anti-concentrated. And instead of inferring that
 471 $p(w_i(y_i - \langle x_i, \ell \rangle)) = 1$, we use that whenever $\pm(y_i - \langle x_i, \ell \rangle) \leq 4\zeta$, $p^2((y_i - \langle x_i, \ell \rangle)) \geq 1 - 4\zeta$.

472 4.1 List-Decodable Regression for Boolean Vectors

473 In this section, we show algorithms for list-decodable regression when the distribution on the
 474 inliers satisfies a weaker anti-concentration condition. This allows us to handle more general inlier
 475 distributions including the product distributions on $\{\pm 1\}^d$, $[0, 1]^d$ and more generally any product
 476 domain. We however require that the unknown linear function be “Boolean”, that is, all its coordinates
 477 be of equal magnitude.

478 We start by defining the weaker anti-concentration inequality. Observe that if $v \in \mathbb{R}^d$ satisfies
 479 $v_i^3 = \frac{1}{d}v_i$ for every i , then the coordinates of v are in $\{0, \pm \frac{1}{\sqrt{d}}\}$.

480 **Definition 4.5** (Certifiable Anti-Concentration for Boolean Vectors). A \mathbb{R}^d valued random variable Y
 481 is k -certifiably (C, δ) -anti-concentrated in *Boolean directions* if there is a univariate polynomial p sat-
 482 isfying $p(0) = 1$ such that there is a degree k sum-of-squares proof of the following two inequalities:
 483 for all $x^2 \leq \delta^2$, $(p(x) - 1)^2 \leq \delta^2$ and for all v such that $v_i^3 = \frac{1}{d}v_i$ for all i , $\|v\|^2 \mathbb{E}_Y p(\langle Y, v \rangle)^2 \leq C\delta$.

484 We can now state the main result of this section.

485 **Theorem 4.6** (List-Decodable Regression in Boolean Directions). *For every α, η , there's a algorithm*
 486 *that takes input a sample generated according to $\text{Lin}_D(\alpha, n, \ell^*)$ in \mathbb{R}^d for D that is k -certifiably*
 487 *$(C, \alpha\eta/10C)$ -anti-concentrated in Boolean directions and $\ell^* \in \left\{ \pm \frac{1}{\sqrt{d}} \right\}^d$ and outputs a list L of size*
 488 *$O(1/\alpha)$ such that there's an $\ell \in L$ satisfying $\|\ell - \ell^*\| < \eta$ with probability at least 0.99 over the*
 489 *draw of the sample. The algorithm requires a sample of size $n \geq (d/\alpha\eta)^{O(\frac{1}{\alpha^2\eta^2})}$ and runs in time*
 490 *$n^{O(k)} = (d/\alpha\eta)^{O(k^2)}$.*

491 The only difference in our algorithm and rounding is that instead of the constraint set $\mathcal{A}_{w,\ell}$, we will
 492 work with $\mathcal{B}_{w,\ell}$ that has an additional constraint $\ell_i^2 = \frac{1}{d}$ for every i . Our algorithm is exactly the
 493 same as Algorithm 1 replacing $\mathcal{A}_{w,\ell}$ by $\mathcal{B}_{w,\ell}$.

$$\mathcal{B}_{w,\ell}: \left\{ \begin{array}{l} \sum_{i=1}^n w_i = \alpha n \\ \forall i \in [n], \quad w_i^2 = w_i \\ \forall i \in [n], \quad w_i \cdot (y_i - \langle x_i, \ell \rangle) = 0 \\ \forall i \in [d], \quad \ell_i^2 = \frac{1}{d} \end{array} \right\} \quad (4.6)$$

494 We will use the following fact in our proof of Theorem 4.6.

495 **Lemma 4.7.** *If a, b satisfy $a^2 = b^2 = \frac{2}{d}$, then, $(a - b)^3 = \frac{1}{d}(a - b)$*

496 *Proof.* $(a - b)^3 = a^3 - b^3 - 3a^2b + 3ab^2 = \frac{1}{d}(a - b - 3b + 3a) = \frac{4}{d}(a - b)$. \square

497 *Proof of Theorem 4.6.* The proof remains the same as in the previous section with one additional
 498 step. First, we can obtain the analog of Lemma 4.1 with a few quick modifications to the proof.
 499 Then, Lemma 4.2 follows from modified Lemma 4.1 as in the previous section. And the proof of
 500 Lemma 4.3 remains exactly the same. We can then put the above lemmas together just as in the proof
 501 of Theorem 1.5.

502 We now describe the modifications to obtain the analog of Lemma 4.1. The key additional step in the
 503 proof of the analog of Lemma 4.1 which follows immediately from Lemma 4.7.

$$\left\{ \forall i \ell_i^2 = \frac{1}{d} \right\} \Big|_{\frac{\ell}{4}} \left\{ (\ell_i - \ell_i^*)^3 = \frac{4}{d}(\ell_i - \ell_i^*) \right\}$$

504 This allows us to replace the usage of certifiable anti-concentration by certifiable anti-concentration
 505 for Boolean vectors and derive:

$$\left\{ \forall i \ell_i^2 = \frac{2}{d} \right\} \Big|_{\frac{\ell}{4}} \left\{ \frac{1}{|\mathcal{I}|} \sum_{i \in \mathcal{I}} p(\langle x_i, \ell - \ell^* \rangle)^2 \leq \frac{\alpha^2 \eta^2}{4} \right\}$$

506 The rest of the proof of Lemma 4.1 remains the same.

507 \square

508 5 Certifiably Anti-Concentrated Distributions

509 In this section, we prove certifiable anti-concentration inequalities for some basic families of distribu-
 510 tions. We first formally state the definition of certified anti-concentration.

511 **Definition 5.1** (Certifiable Anti-Concentration). A \mathbb{R}^d -valued zero-mean random variable Y has a
 512 (C, δ) -anti-concentrated distribution if $\Pr[|\langle Y, v \rangle| \leq \delta \sqrt{\mathbb{E}\langle Y, v \rangle^2}] \leq C\delta$.

513 Y has a k -certifiably (C, δ) -anti-concentrated distribution if there is a univariate polynomial p
 514 satisfying $p(0) = 1$ such that

515 1. $\{ \langle Y, v \rangle^2 \leq \delta^2 \mathbb{E}\langle Y, v \rangle^2 \} \Big|_{\frac{v}{k}} \{ (p(\langle Y, v \rangle) - 1)^2 \leq \delta^2 \}$.

516 2. $\{ \|v\|_2^2 \leq 1 \} \Big|_{\frac{v}{k}} \{ \|v\|_2^2 \mathbb{E} p^2(\langle Y, v \rangle) \leq C\delta \}$.

517 We will say that such a polynomial p “witnesses the certifiable anti-concentration of Y ”. We will
 518 use the phrases “ Y has a certifiably anti-concentrated distribution” and “ Y is a certifiably anti-
 519 concentrated random variable” interchangeably.

520 Before proceeding to prove certifiable anti-concentration of some important families of distributions,
 521 we observe the invariance of the definition under scaling and shifting.

522 **Lemma 5.2** (Scale invariance). *Let Y be a k -certifiably (C, δ) -anti-concentrated random variable.*
 523 *Then, so is cY for any $c \neq 0$.*

524 *Proof.* Let p be the polynomial that witnesses the certifiable anti-concentration of Y . Then, observe
 525 that $q(z) = p(z/c)$ satisfies the requirements of the definition for cY . \square

526 **Lemma 5.3** (Certified anti-concentration of gaussians). *For every $0.1 > \delta > 0$, there is a $k =$
 527 $O\left(\frac{\log^2(1/\delta)}{\delta^2}\right)$ such that $\mathcal{N}(0, I)$ is k -certifiably $(2, 2\delta)$ -anti-concentrated.*

528 *Proof.* Lemma A.1 yields that there exists an univariate even polynomial p of degree k as above such
 529 that for all v , whenever $|\langle x, v \rangle| \leq \delta$, $p(\langle x, v \rangle) \leq 2\delta$, and whenever $\|v\|^2 \leq 1$, $\mathbb{E}_{x \sim \mathcal{N}(0, I)} p(\langle x, v \rangle)^2 \leq$
 530 2δ . Since p is even, $p(z) = \frac{1}{2}(p(z) + p(-z))$ and thus, any monomial in $p(z)$ with non-zero coefficient
 531 must be of even degree. Thus, $p(z) = q(z^2)$ for some polynomial q of degree $k/2$.

532 The first property above for p implies that whenever $z \in [0, \delta]$, $p(z) \leq 2\delta$. By Fact 3.8, we obtain
 533 that:

$$\{\langle x, v \rangle^2 \leq \delta^2\} \Big|_{\frac{v}{k}} \{p(\langle x, v \rangle)^2 \leq \delta\}$$

534 Next, observe that for any j , $\mathbb{E}_{x \sim \mathcal{N}(0, I)} \langle x, v \rangle^{2j} = (2j)!! \cdot \|v\|_2^{2j}$. Thus, $\|v\|_2^2 \mathbb{E}_{x \sim \mathcal{N}(0, I)} p^2(\langle x, v \rangle)$
 535 is a univariate polynomial F in $\|v\|_2^2$. The second property above thus implies that $F(\|v\|_2^2) \leq C\delta$
 536 whenever $\|v\|_2^2 \leq 1$. By another application of Fact 3.8, we obtain:

$$\{\|v\|_2^2 \leq 1\} \Big|_{\frac{v}{k}} \{\mathbb{E}_{x \sim \mathcal{N}(0, I)} p(\langle x, v \rangle)^2 \leq 2\delta\}$$

537 \square

538 We say that Y is a *spherically symmetric* random variable over \mathbb{R}^d if for every orthogonal matrix R ,
 539 RY has the same distribution as Y . Examples include the standard gaussian random variable and
 540 uniform (Haar) distribution on \mathbb{S}^{d-1} . Our argument above for the case of standard gaussian extends
 541 to any distribution that is spherically symmetric and has sufficiently light tails.

542 **Lemma 5.4** (Certified anti-concentration of spherically symmetric, light-tail distributions). *Suppose*
 543 *Y is a \mathbb{R}^d -valued, spherically symmetric random variable such that for any $k \in (0, 2)$, for all t and*
 544 *for all v , $\Pr[\langle v, Y \rangle \geq t\sqrt{\mathbb{E}\langle Y, v \rangle^2}] \leq Ce^{-t^{2/k}/C}$ and for all $\eta > 0$, $\Pr_{x \sim D}[|x| < \eta\sigma] \leq C\eta$, for*
 545 *some absolute constant $C > 0$. Then, for $d = O\left(\frac{\log^{(4+k)/(2-k)}(1/\delta)}{\delta^{2/(2-k)}}\right)$, Y is d -certifiably $(10C, \delta)$ -*
 546 *anti-concentrated.*

547 **Lemma 5.5** (Certified anti-concentration under sampling). *Let D be k -certifiably (C, δ) -anti-*
 548 *concentrated, subexponential and unit covariance distribution. Let S be a collection of n independent*
 549 *samples from D . Then, for $n \geq \Omega((kd \log(d))^{O(k)})$, with probability at least $1 - 1/d$, the uniform*
 550 *distribution on S is $(2C, \delta)$ -anti-concentrated.*

551 *Proof.* Let p be the degree k polynomial that witnesses the certifiable anti-concentration of D . Let Y
 552 be the random variable with distribution D' , the uniform distribution on n i.i.d. samples from D . We
 553 will show that p also witnesses that k -certifiable $(4C, \delta/2)$ -anti-concentration of Y . To this end it is
 554 sufficient to take enough samples such that the following holds.

$$\Pr\left(|\mathbb{E}_D[p^2(\langle Y, v \rangle)] - \mathbb{E}_{D'}[p^2(\langle Y, v \rangle)]| > \mathbb{E}_D[p^2(\langle Y, v \rangle)]/2\right) < 1/d$$

555 Observe that $p^2(\langle Y, v \rangle)$ may be written as $\langle c(Y)c(Y)^T, m(v)m(v)^T \rangle$ where $c(Y)$ are the coefficients
 556 of $p(\langle Y, v \rangle)$ and $m(v)$ is the vector containing monomials. The dot product above is the usual trace
 557 inner product between matrices. Now, it is sufficient to show that

$$\Pr\left(\|\mathbb{E}_{D'}c(Y)c(Y)^T - \mathbb{E}_Dc(Y)c(Y)^T\|_F^2 > \|\mathbb{E}_Dc(Y)c(Y)^T\|_F^2/4\right) < 1/d$$

558 Since p was a univariate polynomial of degree k in d dimensional variables, there are at most d^{2k}
 559 entries in total, and each entry is at most a degree $2k$ polynomial of subexponential random variables
 560 in d variables. Using standard concentration results for polynomials of subexponential random
 561 variables (for instance Theorem 1.2 from [27] and the references therein). We see that each entry
 562 satisfies

$$\Pr\left(|\mathbb{E}_Dc(Y)_i c(Y)_j - \mathbb{E}_{D'}c(Y)_i c(Y)_j| > \epsilon\right) \leq \exp\left(-\Omega\left(\frac{n\epsilon}{\mathbb{E}(c(Y)_i c(Y)_j)^2}\right)^{1/2k}\right)$$

563 An application of a union bound, squaring the term inside and replacing ϵ^2 by $\mathbb{E}(c(Y)_i c(Y)_j)^2/4$
 564 gives us

$$\Pr \left(\sum_{i,j=1}^{d^{2k}} (\mathbb{E}_D c(Y)_i c(Y)_j - \mathbb{E}_{D'} c(Y)_i c(Y)_j)^2 > \|\mathbb{E} c(Y) c(Y)^T\|_F^2/4 \right) \leq d^{2k} \exp \left(-\Omega \left(\frac{n}{d^{O(k)}} \right)^{1/2k} \right)$$

565 Hence, setting $n = O((kd \log(d))^{O(k)})$ ensures that with probability at least $1 - 1/d$, the distribution
 566 D' is $(2C, \delta)$ -anti-concentrated.

567 □

568 We say that a $d \times d$ matrix A is C' -well-conditioned if all singular values of A are within a factor of
 569 C' of each other.

570 **Lemma 5.6** (Certified anti-concentration under linear transformations). *Let Y be k -certifiably (C, δ) -*
 571 *anti-concentrated random variable over \mathbb{R}^d . Let A be any C' -well-conditioned linear transformation.*
 572 *Then, AY is k -certifiably $(C, C'^2\delta)$ -anti-concentrated.*

573 *Proof.* Let $\|A\|$ be the largest singular value of A . Let p be a polynomial that witnesses the certifiable
 574 anti-concentration of Y . Let $q(z) = p(z/\|A\|)$. We will prove that q witnesses the k -certifiable
 575 $(C, C'^2\delta)$ -anti-concentration of AY .

576 Towards this, observe that:

$$\{ \langle Y, v \rangle^2 \leq \delta^2 \mathbb{E} \langle Y, v \rangle^2 \} \Big|_{\frac{v}{2}} \{ \langle AY, v \rangle^2 \leq \delta^2 \mathbb{E} \langle AY, v \rangle^2 \} .$$

577

$$\{ \langle Y, (A^T v)/\|A\| \rangle^2 \leq \delta^2 \mathbb{E} \langle Y, (A^T v)/\|A\| \rangle^2 \} \Big|_{\frac{v}{k}} \{ (p(\langle Y, (A^T v)/\|A\| \rangle) - 1)^2 \leq \delta^2 \} ,$$

578 this is the same as

$$\{ \langle AY, v \rangle^2 \leq \delta^2 \mathbb{E} \langle AY, v \rangle^2 \} \Big|_{\frac{v}{k}} \{ (q(\langle AY, v \rangle) - 1)^2 \leq \delta^2 \} .$$

579 Where $q = p(x/\|A\|)$. Now, for $w = (A^T v)/\|A\|$ and any unit vector v ,

$$\{ \|w\|_2^2 \leq 1 \} \Big|_{\frac{v}{k}} \{ \|A^T v\|_2^2 / \|A\|_2^2 \mathbb{E} p^2(\langle AY, v \rangle / \|A\|) \leq C\delta \} ,$$

580 Thus,

$$\{ \|A^T v\|_2^2 \leq \|A\|^2 \} \Big|_{\frac{v}{k}} \{ \|A^T v\|_2^2 \mathbb{E} q^2(\langle AY, v \rangle) \leq C \|A\|_2^2 \delta \} .$$

581 However,

$$\{ \|v\|_2^2 \leq 1 \} \Big|_{\frac{v}{2}} \{ \|A^T v\|_2^2 \leq \|A\|^2 \} ,$$

582 and thus,

$$\{ \|v\|_2^2 \leq 1 \} \Big|_{\frac{v}{k}} \{ \|v\|_2^2 \mathbb{E} q^2(\langle AY, v \rangle) \leq CC'^2\delta \} .$$

583 □

584 **Lemma 5.7** (Certifiable Anti-Concentration in Boolean Directions). *Fix $C > 0$. Let Y be a \mathbb{R}^d*
 585 *valued product random variable satisfying:*

- 586 1. **Identical Coordinates:** Y_i are identically distributed for every $1 \leq i \leq d$.
- 587 2. **Anti-Concentration** For every $v \in \left\{ 0, \pm \frac{1}{\sqrt{d}} \right\}^d$, $\Pr[|\langle Y, v \rangle| \leq \delta \sqrt{\mathbb{E} \langle Y, v \rangle^2}] \leq C\delta$.
- 588 3. **Light tails** For every $v \in \mathbb{S}^{d-1}$, $\Pr[|\langle Y, v \rangle| > t \sqrt{\mathbb{E} \langle Y, v \rangle^2}] \leq \exp(-t^2/C)$.

589 Then, Y is k -certifiably (C, δ) -anti-concentrated for $k = O\left(\frac{\log^2(1/\delta)}{\delta^2}\right)$.

590 *Proof.* We use the p from Lemma A.1. To see that p witnesses the anti-concentration of Y , once
591 again observe that Lemma A.1 applies to give us a real life proof of the required statements. We
592 now exhibit a sum of squares proof. Observe that every monomial of even degree $2k$ for any
593 $k \in \mathbb{N}$, $\mathbb{E}_{Y \sim D} \langle Y, v \rangle^{2k}$ is a *symmetric* polynomial in v with non-zero coefficients only on even-degree
594 monomials in v . This follows by noting that the coordinates of D are independent and identically
595 distributed and x^2 is an even function. It is a fact that all symmetric polynomials in v can be expressed
596 as polynomials in the “power-sum” polynomials $\|v\|_2^{2i}$ for $i \leq 2t$. However, since $v_i^2 \in \{0, \frac{1}{d}\}$ for
597 $i \geq 1$, $\|v\|_2^{2i} = \frac{1}{d^{i-1}} \|v\|_2^2$. Hence a polynomial in $\|v\|_2^{2i}$ is also a univariate polynomial in $\|v\|_2^2$.
598 Since these are polynomial inequalities, they are also sum-of-squares proofs of these inequalities.

599 The observation above implies $\|v\|_2^2 \mathbb{E}_Y p(\langle Y, v \rangle)^2 = \|v\|_2^2 \cdot F(\|v\|_2^2)$ for some degree k univariate
600 polynomial F . Since F is a univariate polynomial and $\|v\|_2^2 \leq 1$ is an “interval constraint” by
601 applying Fact 3.8, we get: $\frac{\|v\|_2^2}{2t} \{ \|v\|_2^2 F(\|v\|_2^2) \leq C\delta \}$. Recalling the fact that $\|v\|_2^2 \mathbb{E}_Y p(\langle Y, v \rangle)^2 =$
602 $\|v\|_2^2 \cdot F(\|v\|_2^2)$, this completes the proof. \square

603 6 Information-Theoretic Lower Bounds for List-Decodable Regression

604 In this section, we show that list-decodable regression on $\text{Lin}_D(\alpha, \ell^*)$ information-theoretically
605 requires that D satisfy α -anti-concentration: $\Pr_{x \sim D}[\langle x, v \rangle = 0] < \alpha$ for any non-zero v .

606 **Theorem 6.1 (Main Lower Bound).** *For every q , there is a distribution D on \mathbb{R}^d satisfying*
607 $\Pr_{x \sim D}[\langle x, v \rangle = 0] \leq \frac{1}{q}$ *such that there’s no $\frac{1}{2q}$ -approximate list-decodable regression algorithm for*
608 $\text{Lin}_D(\frac{1}{q}, \ell^*)$ *that can output a list of size $< d$.*

609 *Remark 6.2 (Impossibility of Mixed Linear Regression on the Hypercube).* Our construction for the
610 case of $q = 2$ actually shows the impossibility of the well-studied and potentially easier problem of
611 noiseless *mixed linear regression* on the uniform distribution on $\{0, 1\}^n$. This is because \mathcal{R}_i is, by
612 construction, obtained by using one of e_i or $1 - e_i$ to label each example point with equal probability.

613 Theorem 6.1 is tight in a precise way. In Proposition 2.4, we proved that whenever D satisfies
614 $\Pr_{x \sim D}[\langle x, v \rangle = 0] < \frac{1}{q}$, there is an (inefficient) algorithm for *exact* list-decodable regression
615 algorithm for $\text{Lin}_D(\frac{1}{q}, \ell^*)$. Note that our lower bound holds even in the setting where there is no
616 additive noise in the inliers.

617 Somewhat surprisingly, our lower bound holds for extremely natural and well-studied distributions -
618 uniform distribution on $\{0, 1\}^n$ and more generally, uniform distribution on $\{0, 1, \dots, q-1\}^d = [q]^d$
619 for any q . We can easily determine a tight bound on the anti-concentration of both these distributions.

620 **Lemma 6.3.** *For any non-zero $v \in \mathbb{R}^d$, $\Pr_{x \sim \{0, 1\}^n} \langle x, v \rangle = 0 \leq \frac{1}{2}$ and $\Pr_{x \sim [q]^d} [\langle x, v \rangle = 0] \leq \frac{1}{q}$.*

621 Note that this is tight for any $v = e_i$, the vector with 1 in the i th coordinates and 0s in all others.

622 *Proof.* Fix any v . Without loss of generality, assume that all coordinates of v are non-zero. If not, we
623 can simply work with the uniform distribution on the sub-hypercube corresponding to the non-zero
624 coordinates of v .

625 Let $S \subseteq \{0, 1\}^n$ ($[q]^d$, respectively) be the set of all $x \in \{0, 1\}^n$ ($[q]^d$, respectively) such that
626 $\langle x, v \rangle = 0$. Then, observe that for any $x \in S$, and any i , $x^{(i)}$ obtained by flipping the i th bit
627 (changing the i th coordinate to any other value) of x cannot be in S . Thus, S is an independent set in
628 the graph on $\{0, 1\}^n$ (in $[q]^d$, respectively) with edges between pairs of points with hamming distance
629 1.

630 It is a standard fact [56] that the maximum independent set in the d -hypercube is of size exactly
631 2^{d-1} and in the q -ary Hamming graph $[q]^d$ is of size q^{d-1} . Thus, $\Pr_{x \sim \{0, 1\}^n} [\langle x, v \rangle = 0] \leq \frac{1}{2}$ and
632 $\Pr_{x \sim [q]^d} [\langle x, v \rangle = 0] \leq \frac{1}{q}$.

633 \square

634 To prove our lower bound, we give a family of d distributions on labeled linear equations, \mathcal{R}_i for
635 $1 \leq i \leq d$ that satisfy the following:

- 636 1. The examples in each are chosen from uniform distribution on $[q]^d$,
637 2. $\frac{1}{q}$ fraction of the samples are labeled by e_i in \mathcal{R}_i , and,
638 3. for any i, j , \mathcal{R}_i and \mathcal{R}_j are statistically indistinguishable.

639 Thus, given samples from \mathcal{R}_i , any $\frac{1}{2q}$ -approximate list-decoding algorithm must produce a list of
640 size at least d .

641 Our construction and analysis of \mathcal{R}_i is simple and exactly the same in both the cases. However
642 it is somewhat easier to understand for the case of the hypercube ($q = 2$). The following simple
643 observation is the key to our construction.

644 **Lemma 6.4.** *For $1 \leq i \leq d$, let \mathcal{R}_i be the distribution on linear equations induced by the following
645 sampling method: Sample $x \sim \{0, 1\}^d$, choose $a \sim \{0, 1\}$ uniformly at random and output:
646 $(x, \langle x, (1 - a)e_i \rangle)$. Then, $\mathcal{R}_i = \mathcal{R}_j$ for any $i, j \leq d$.*

647 *Proof.* The proof follows by observing that \mathcal{R}_i when viewed as a distribution on \mathbb{R}^{d+1} is same as the
648 uniform distribution on $\{0, 1\}^{d+1}$ and thus independent of i . \square

649 The argument immediately generalizes to $[q]^d$ and yields:

650 **Lemma 6.5.** *For $1 \leq i \leq d$, let \mathcal{R}_i be the distribution on linear equations induced by the fol-
651 lowing sampling method: Sample $x \sim [q]^d$, choose $a \sim \{0, 1\}$ uniformly at random and output:
652 $(x, (\langle x, e_i \rangle + a) \bmod q)$. Then, $\mathcal{R}_i = \mathcal{R}_j$ for any $i, j \leq d$.*

653 In this case, we interpret the $1/q$ fraction of the samples where $a = 0$ as the inliers. Observe that these
654 are labeled by a single linear function e_i in any \mathcal{R}_i . Thus, they form a valid model in $\text{Lin}_D(\alpha, \ell^*)$ for
655 $\alpha = 1/q$.

656 Since the linear functions defined by e_i on $[q]^d$, when normalized to have unit norm, have a pairwise
657 Euclidean distance of at least $1/q$, we immediately obtain a proof of Theorem 6.1.

658 A Polynomial Approximation for Core-Indicator

659 The main result of this section is a low-degree polynomial approximator for the function $\mathbf{1}(|x| < \delta)$
660 with respect to all distributions that have asymptotically lighter-than-exponential tails.

661 **Lemma A.1.** *Let D be a distribution on \mathbb{R} with mean 0, variance $\sigma^2 \leq 1$ and satisfying:*

662 1. **Anti-Concentration:** For all $\eta > 0$, $\Pr_{x \sim D}[|x| < \eta\sigma] \leq C\eta$, and,

663 2. **Tail bound:** $\Pr[|x| \geq t\sigma] \leq e^{-\frac{t^2/k}{C}}$ for $k < 2$ and all t ,

664 for some $C > 1$. Then, for any $\delta > 0$, there is a $d = O\left(\frac{\log^{(4+k)/(2-k)}(1/\delta)}{\delta^{2/(2-k)}}\right) = \tilde{O}\left(\frac{1}{\delta^{2/(2-k)}}\right)$
665 and an even polynomial $q(x)$ of degree d such that $q(0) = 1$, $q(x) = 1 \pm \delta$ for all $|x| \leq \delta$ and
666 $\sigma^2 \cdot \mathbb{E}_{x \sim D}[q^2(x)] \leq 10C\delta$.

667 Before proceeding to the proof, we note that the bounds on the degree above are tight up to poly
668 logarithmic factors for the gaussian distribution.

669 **Lemma A.2.** *For every polynomial p of degree d such that $p(0) = 1$, $\mathbb{E}_{x \sim \mathcal{N}(0,1)}[p^2(x)] = \Omega\left(\frac{1}{\sqrt{d}}\right)$.
670 Further, there is a polynomial p_* of degree d such that $p_*(0) = 1$ and $\mathbb{E}_{x \sim \mathcal{N}(0,1)}p_*^2(x) = \Theta\left(\frac{1}{\sqrt{d}}\right)$.*

671 Our construction of the polynomial is based on standard techniques in approximation theory for
672 constructing polynomial approximators for continuous functions over an interval. Most relevant for
673 us are various works of Eremenko and Yuditskii [24, 25, 23] and Diakonikolas, Gopalan, Jaiswal,
674 Servedio and Viola [13] on such constructions for the sign function on the interval $[-1, a] \cup [a, 1]$ for
675 $a > 0$. We point the reader to the excellent survey of this beautiful line of work by Lubinsky [43].

676 **Fact A.3** (Theorem 3.5 in [13]). Let $0 < \eta < 0.1$, then there exist constants C, c such that for

$$a := \eta^2 / C \log(1/\eta) \text{ and } K = 4c \log(1/\eta) / a + 2 < O(\log^2(1/\eta) / \eta^2)$$

677 there is a polynomial $p(t)$ of degree K satisfying

- 678 1. $p(t) > \text{sign}(t) > -p(-t)$ for all $t \in \mathbb{R}$.
- 679 2. $p(t) \in [\text{sign}(t), \text{sign}(t) + \eta]$ for $t \in [-1/2, -2a] \cup [0, 1/2]$.
- 680 3. $p(t) \in [-1, 1 + \eta]$ for $t \in (-2a, 0)$
- 681 4. $|p(t)| \leq 2 \cdot (4t)^K$ for all $t > \frac{1}{2}$.

682 We will also rely on the following elementary integral estimate.

Lemma A.4 (Tail Integral).

$$\int_{[L, \infty]} \exp\left(-\frac{x^{2/k}}{C}\right) x^{2d} dx < \exp\left(-\frac{L^{2/k}}{C}\right) ((L)^{4d} + (16kd)^{kd}).$$

683 *Proof.* We first prove the claim for $k = 1$. Let $y = x - L$. The, $\int_L^\infty e^{-x^2} x^{2d} dx = \int_0^\infty e^{-(y+L)^2} (y +$
684 $L)^{2d} dy$. We now use that $y^2 + L^2 \leq (y + L)^2$ for all $y \geq 0$ and $(y + L)^{2d} \leq 2^{2d}(y^{2d} + L^{2d})$ to
685 upper bound the integral above by: $e^{-L^2} L^{2d} + 2^{2d} e^{-L^2} \int_0^\infty e^{-y^2} y^{2d}$. Using $\int_0^\infty e^{-y^2} y^{2d} < (4d)^d$
686 gives a bound of $e^{-L^2} (L^{2d} + (8d)^d)$.

687 For larger k , we substitute $y = x^{1/k}$ and write the integral in question as $\int_{L^{1/k}}^\infty e^{-y^2} y^{2kd - (k-1)} dy$.
688 Applying the calculation from the above special case, this integral is upper bounded by: $e^{-L^{2/k}} (L^{4d} +$
689 $(16kd)^{kd})$. \square

690 *Proof of Lemma A.1.* Let $p(x)$ be the degree $d < O\left(\frac{L \log^2(1/\delta)}{\delta}\right)$ polynomial from Fact A.3. We
691 then construct a polynomial $q(x)$ that will be close to 0 in the range $[\delta, L]$ and $[-L, -\delta]$ and close to
692 1 in the range $[-\delta, \delta]$. Our polynomial q is obtained by shifting and appropriately scaling two copies
693 of p .

$$q(x) = \frac{p\left(a + \frac{x}{4L}\right) + p\left(-\left(a + \frac{x}{4L}\right)\right) - 1}{p(a) + p(-a) - 1}$$

694 Then, $q(0) = 1$. It further satisfies:

- 695 1. $q(x) \in [0, C\sqrt{\delta/L}]$ for $x \in [\delta, L] \cup [-L, \delta]$.
- 696 2. $q(x) \in [1 - C\sqrt{\delta/L}, 1 + \sqrt{\delta/L}]$ for $x \in [-\delta, \delta]$.
- 697 3. $q(x) \in [0, 1 + \sqrt{\delta/L}]$ for $x \in [-3\delta, -\delta] \cup [\delta, 3\delta]$.
- 698 4. $|q(x)| < 4 \cdot (4x)^t$ for $|x| > L$

699 We now prove the bound the $\mathbb{E}p^2$. We do this by providing upper bounds on the contributions to
700 $\sigma^2 \cdot \mathbb{E}_{x \sim \mathcal{D}} [q^2(\sigma x)]$ from the disjoint sets with different guarantees below. Since we are going to
701 evaluate $q(\sigma x)$ the intervals will be scaled by σ .

702 The contributions from the regions $\frac{1}{\sigma}[\delta, L]$ and $\frac{1}{\sigma}[-\delta, \delta]$ can be naively upper bounded by the
703 maximum value that the polynomial can take here times the probability of landing in these regions.
704 The first of these contributes $\sigma \cdot \frac{\delta}{L} \cdot (L - \delta) \leq \delta$, and using anticoncentration, the second region
705 contributes $\left(1 + \sqrt{\frac{\delta}{L}}\right)^2 \cdot 2C\delta \leq 4C\delta$. The region $\frac{1}{\sigma}[\delta, 3\delta]$ can be bounded similarly to get an upper
706 bound of $2 \left(1 + \sqrt{\frac{\delta}{L}}\right)^2 \sigma^2 \delta \leq 4\delta$. To finish, we use Lemma A.4 to upper bound the contribution to

707 $\mathbb{E}p^2$ from the tail:

$$\begin{aligned} \sigma^2 C' \int_{\frac{1}{\sigma}[L, \infty]} q^2(\sigma x) \exp\left(-\frac{x^{2/k}}{C}\right) dx &\lesssim \sigma^{2+d} 4^d \exp\left(-\frac{1}{C} \cdot \left(\frac{L}{\sigma}\right)^{2/k}\right) ((L/\sigma)^{4d} + (16kd)^{kd}) \\ &\lesssim \exp\left(2d + 4d \log\left(\frac{L}{\sigma}\right) - \frac{1}{C} \cdot \left(\frac{L}{\sigma}\right)^{2/k} + kd \log(16kd)\right) \end{aligned}$$

708 We choose L satisfying $10d \log(d) + 4d \log(\frac{L}{\sigma}) - \frac{1}{C} \cdot (\frac{L}{\sigma})^{2/k} < 2 \log(1/\delta)$.

709 Since $d = O\left(\frac{L \log^2(1/\delta)}{\delta}\right)$, $k < 2$, and $\sigma < 1$ we can now choose $L = \left(\frac{C100 \log^3(1/\delta)}{\delta}\right)^{k/(2-k)}$ to

710 satisfy the inequality above and to get $d \lesssim \frac{\log^{2+3k/(2-k)}(1/\delta)}{\delta^{1+k/(2-k)}}$. When $k = 1$ we get $d = \tilde{O}(1/\delta^2)$.

711 Since $\sigma < 1$ in all the above calculations, we get our result by re-scaling δ .

712 \square

713 We now complete the proof of Lemma A.2.

714 *Proof of Lemma A.2.* Any polynomial p of degree d can be written as $p(x) = \sum_{i=1}^d \alpha_i h_i(x)$ where
715 h_i denote the hermite polynomials of degree i , satisfying $\mathbb{E}_{x \sim \mathcal{N}(0,1)} h_i = 0$ and $\mathbb{E}_{x \sim \mathcal{N}(0,1)} [h_i^2(x)] =$
716 1. Since $p(0) = 1$, using Cauchy-Schwartz inequality, we obtain:

$$\mathbb{E}_{x \sim \mathcal{N}(0,1)} [p^2(x)] \cdot \sum_{i=1}^d h_i^2(0) = \left(\sum_{i=1}^d \alpha_i\right) \cdot \left(\sum_{i=1}^d h_i^2(0)\right) \geq \left(\sum_{i=1}^d \alpha_i h_i(0)\right)^2 \geq 1$$

717 Further, observe that for the polynomial $p_*(x) = \frac{1}{\sum_{i=1}^d h_i^2(0)} \sum_{i=1}^d h_i(0) h_i(x)$, the above inequality is

718 tight. Using that $h_{2i}(0) = \frac{(2i-1)!!}{\sqrt{(2i)}}$ and $h_i(0) = 0$ if i is odd, (see, for e.g., [55]), we have:

$$\begin{aligned} \mathbb{E}_{x \sim \mathcal{N}(0,1)} [p^2(x)] &\geq \mathbb{E}_{x \sim \mathcal{N}(0,1)} p_*^2(x) = \left(\sum_{i=1}^d h_i^2(0)\right)^{-1} = \left(\sum_{i=1}^{d/2} \left(\frac{(2i-1)!!}{\sqrt{(2i)}}\right)^2\right)^{-1} \\ &= \left(\sum_{i=1}^{d/2} \frac{(2i)!}{2^{2i} i!^2}\right)^{-1} = \left(\sum_{i=1}^{d/2} \binom{2i}{i} \cdot \frac{1}{2^{2i}}\right)^{-1} = \Theta\left(\sum_{i=1}^{d/2} \frac{1}{\sqrt{i}}\right)^{-1} = \Theta(\sqrt{d})^{-1}. \end{aligned}$$

719 \square

720 B Brute-force search can generate a $\exp(d)$ size list

721 In the following, we write e_i to denote the vector with 1 in the i th coordinate and 0s in all others.

722 **Proposition B.1.** *There exists a distribution D on \mathbb{R}^d and a model $\text{Lin}_D(\alpha, \ell^*)$ such that for every*
723 $\alpha < 1/2$, *with probability at least $1 - 1/d$ over the draw of a n -size sample \mathcal{S} from $\text{Lin}_D(\alpha, \ell^*)$, there*
724 *exists a collection $\text{Sol} \subseteq \{S \subseteq \mathcal{S} \mid |S| = \alpha n\}$ of size $\exp(d)$ and unit length vectors ℓ_S for every*
725 $S \in \text{Sol}$ *such that ℓ_S satisfies all equations in S and for every $S \neq S' \in \text{Sol}$, $\|\ell_S - \ell_{S'}\|_2 \geq 0.1$.*

726 *Proof.* Let D be the uniform distribution on $e_1, e_2, \dots, e_d \in \mathbb{R}^d$. Let $\ell^* := \vec{1}/\sqrt{d}$ be the all-ones
727 vector in \mathbb{R}^d scaled by $1/\sqrt{d}$ and let d samples be drawn from the uncorrupted distribution. These give
728 us our inliers, $\mathcal{I} = \{(x_i, y_i)\}_{i=1}^{\alpha n}$. For the outliers, choose the following multiset $\mathcal{O} := 1/\alpha - 1$ copies
729 of $\{(e_i, j) \mid i \in [d], j \in \{\pm 1/\sqrt{d}\}\}$. This is a sample set of size $2d/\alpha$. Any $a \in \{\pm 1/\sqrt{d}\}^d$ is a valid
730 candidate for a solution for this data. This is because for any such a , $\mathcal{I}_a := \{(e_i, a_i) \mid i \in [d]\} \subset \mathcal{S}$
731 satisfies the following

732 1. $\mathcal{I}_a \subset \mathcal{S}$, $|\mathcal{I}_a| = d = \frac{\alpha}{2} |\mathcal{S}|$ and

733 2. for any $(x, y) \in \mathcal{I}_a$, $y = \langle x, a \rangle$.

734 The Gilbert–Varshamov bound from coding theory now tells us that there are at least $\Omega(\exp(\Omega(d)))$
735 $\{0, 1\}$ vectors in d dimensions that pairwise have a hamming distance of $0.1 \cdot d$. This transfers to the
736 set $\{\pm 1/\sqrt{d}\}$ to give us that there are $\Omega(\exp(\Omega(d)))$ vectors in $\{\pm 1/\sqrt{d}\}$ that are pairwise 0.1 apart
737 in 2-norm.

738

□

739 References

- 740 [1] Pranjali Awasthi, Maria-Florina Balcan, and Philip M. Long. The power of localization for
741 efficiently learning linear separators with malicious noise. *CoRR*, abs/1307.8371, 2013. 1
- 742 [2] Sivaraman Balakrishnan, Martin J. Wainwright, and Bin Yu. Statistical guarantees for the EM
743 algorithm: From population to sample-based analysis. *CoRR*, abs/1408.2156, 2014. 2
- 744 [3] Maria-Florina Balcan, Avrim Blum, and Santosh Vempala. A discriminative framework for
745 clustering via similarity functions. In *STOC*, pages 671–680. ACM, 2008. 1
- 746 [4] Boaz Barak, Jonathan A. Kelner, and David Steurer. Dictionary learning and tensor decomposi-
747 tion via the sum-of-squares method [extended abstract]. In *STOC’15—Proceedings of the 2015*
748 *ACM Symposium on Theory of Computing*, pages 143–151. ACM, New York, 2015. 5
- 749 [5] Boaz Barak and Ankur Moitra. Noisy tensor completion via the sum-of-squares hierarchy. In
750 *COLT*, volume 49 of *JMLR Workshop and Conference Proceedings*, pages 417–445. JMLR.org,
751 2016. 5
- 752 [6] Boaz Barak and David Steurer. Proofs, beliefs, and algorithms through the lens of sum-of-
753 squares, 2016. Lecture notes in preparation, available on <http://sumofsquares.org>.
754 8
- 755 [7] Thorsten Bernholt. Robust estimators are hard to compute. Technical report, Technical
756 Report/Universität Dortmund, SFB 475 Komplexitätsreduktion in Multivariaten Datenstrukturen,
757 2006. 1
- 758 [8] Kush Bhatia, Prateek Jain, Parameswaran Kamalaruban, and Purushottam Kar. Consistent robust
759 regression. In *Advances in Neural Information Processing Systems 30: Annual Conference*
760 *on Neural Information Processing Systems 2017, 4-9 December 2017, Long Beach, CA, USA*,
761 pages 2107–2116, 2017. 2
- 762 [9] Moses Charikar, Jacob Steinhardt, and Gregory Valiant. Learning from untrusted data. In *STOC*,
763 pages 47–60. ACM, 2017. 1, 2
- 764 [10] Yudong Chen, Xinyang Yi, and Constantine Caramanis. A convex formulation for mixed
765 regression with two components: Minimax optimal rates. In *Proceedings of The 27th Conference*
766 *on Learning Theory, COLT 2014, Barcelona, Spain, June 13-15, 2014*, pages 560–604, 2014. 2
- 767 [11] Yu Cheng, Ilias Diakonikolas, and Rong Ge. High-dimensional robust mean estimation in
768 nearly-linear time. In *Proceedings of the Thirtieth Annual ACM-SIAM Symposium on Discrete*
769 *Algorithms, SODA 2019, San Diego, California, USA, January 6-9, 2019*, pages 2755–2771,
770 2019. 1
- 771 [12] Richard D. De Veaux. Mixtures of linear regressions. *Comput. Statist. Data Anal.*, 8(3):227–245,
772 1989. 2
- 773 [13] Ilias Diakonikolas, Parikshit Gopalan, Ragesh Jaiswal, Rocco A. Servedio, and Emanuele Viola.
774 Bounded independence fools halfspaces. In *50th Annual IEEE Symposium on Foundations of*
775 *Computer Science, FOCS 2009, October 25-27, 2009, Atlanta, Georgia, USA*, pages 171–180,
776 2009. 18, 19
- 777 [14] Ilias Diakonikolas, Gautam Kamath, Daniel M. Kane, Jerry Li 0001, Jacob Steinhardt, and Alis-
778 tair Stewart. Sever: A robust meta-algorithm for stochastic optimization. *CoRR*, abs/1803.02815,
779 2018. 2

- 780 [15] Ilias Diakonikolas, Gautam Kamath, Daniel M. Kane, Jerry Li, Ankur Moitra, and Alistair
781 Stewart. Robust estimators in high dimensions without the computational intractability. In
782 *FOCS*, pages 655–664. IEEE Computer Society, 2016. [1](#)
- 783 [16] Ilias Diakonikolas, Gautam Kamath, Daniel M. Kane, Jerry Li, Ankur Moitra, and Alistair
784 Stewart. Robustly learning a gaussian: Getting optimal error, efficiently. *CoRR*, abs/1704.03866,
785 2017. [1](#)
- 786 [17] Ilias Diakonikolas, Gautam Kamath, Daniel M. Kane, Jerry Li, Ankur Moitra, and Alistair
787 Stewart. Robustly learning a gaussian: Getting optimal error, efficiently. In *Proceedings of*
788 *the Twenty-Ninth Annual ACM-SIAM Symposium on Discrete Algorithms, SODA 2018, New*
789 *Orleans, LA, USA, January 7-10, 2018*, pages 2683–2702, 2018. [1](#)
- 790 [18] Ilias Diakonikolas, Gautam Kamath, Daniel M. Kane, Jerry Zheng Li, Ankur Moitra, and
791 Alistair Stewart. Robust estimators in high dimensions without the computational intractability.
792 *CoRR*, abs/1604.06443, 2016. [1](#)
- 793 [19] Ilias Diakonikolas, Daniel M. Kane, and Alistair Stewart. Learning geometric concepts with
794 nasty noise. *CoRR*, abs/1707.01242, 2017. [1](#)
- 795 [20] Ilias Diakonikolas, Daniel M. Kane, and Alistair Stewart. List-decodable robust mean estimation
796 and learning mixtures of spherical gaussians. In *Proceedings of the 50th Annual ACM SIGACT*
797 *Symposium on Theory of Computing, STOC 2018, Los Angeles, CA, USA, June 25-29, 2018*,
798 pages 1047–1060, 2018. [1](#), [2](#)
- 799 [21] Ilias Diakonikolas, Weihao Kong, and Alistair Stewart. Efficient algorithms and lower bounds
800 for robust linear regression. In Timothy M. Chan, editor, *Proceedings of the Thirtieth Annual*
801 *ACM-SIAM Symposium on Discrete Algorithms, SODA 2019, San Diego, California, USA,*
802 *January 6-9, 2019*, pages 2745–2754. SIAM, 2019. [2](#)
- 803 [22] P. Erdős. On a lemma of littlewood and offord. *Bull. Amer. Math. Soc.*, 51(12):898–902, 12
804 1945. [3](#), [4](#)
- 805 [23] Alexandre Eremenko and Peter Yuditskii. Uniform approximation of $\operatorname{sgn} x$ by polynomials and
806 entire functions. *J. Anal. Math.*, 101:313–324, 2007. [18](#)
- 807 [24] Alexandre Eremenko and Peter Yuditskii. An extremal problem for a class of entire functions.
808 *C. R. Math. Acad. Sci. Paris*, 346(15-16):825–828, 2008. [18](#)
- 809 [25] Alexandre Eremenko and Peter Yuditskii. Polynomials of the best uniform approximation to
810 $\operatorname{sgn}(x)$ on two intervals. *J. Anal. Math.*, 114:285–315, 2011. [18](#)
- 811 [26] Susana Faria and Gilda Soromenho. Fitting mixtures of linear regressions. *J. Stat. Comput.*
812 *Simul.*, 80(1-2):201–225, 2010. [2](#)
- 813 [27] Friedrich Götze, Holger Sambale, and Arthur Sinulis. Concentration inequalities for polynomials
814 in α -sub-exponential random variables. *arXiv e-prints*, page arXiv:1903.05964, Mar 2019. [15](#)
- 815 [28] M. Grötschel, L. Lovász, and A. Schrijver. The ellipsoid method and its consequences in
816 combinatorial optimization. *Combinatorica*, 1(2):169–197, 1981. [8](#), [9](#)
- 817 [29] Frank R Hampel, Elvezio M Ronchetti, Peter J Rousseeuw, and Werner A Stahel. *Robust*
818 *statistics: the approach based on influence functions*, volume 114. John Wiley & Sons, 2011. [1](#)
- 819 [30] Sam B. Hopkins and Jerry Li. Mixture models, robustness, and sum of squares proofs. 2017. [1](#),
820 [2](#), [5](#)
- 821 [31] Peter J Huber. Robust statistics. In *International Encyclopedia of Statistical Science*, pages
822 1248–1251. Springer, 2011. [1](#)
- 823 [32] Michael I. Jordan and Robert A. Jacobs. Hierarchical mixtures of experts and the em algorithm.
824 *Neural Computation*, 6(2):181–214, 1994. [2](#)

- 825 [33] Sushrut Karmalkar and Eric Price. Compressed sensing with adversarial sparse noise via l1
826 regression. In Jeremy T. Fineman and Michael Mitzenmacher, editors, *SOSA@SODA*, volume 69
827 of *OASICS*, pages 19:1–19:19. Schloss Dagstuhl - Leibniz-Zentrum fuer Informatik, 2019. [2](#)
- 828 [34] Adam R. Klivans, Pravesh K. Kothari, and Raghu Meka. Efficient algorithms for outlier-robust
829 regression. In *Conference On Learning Theory, COLT 2018, Stockholm, Sweden, 6-9 July 2018.*,
830 pages 1420–1430, 2018. [1](#), [2](#), [4](#), [5](#)
- 831 [35] Adam R. Klivans, Philip M. Long, and Rocco A. Servedio. Learning halfspaces with malicious
832 noise. *Journal of Machine Learning Research*, 10:2715–2740, 2009. [1](#)
- 833 [36] Pravesh K. Kothari and Jacob Steinhardt. Better agnostic clustering via relaxed tensor norms.
834 2017. [1](#), [2](#), [4](#), [5](#)
- 835 [37] Pravesh K. Kothari and David Steurer. Outlier-robust moment-estimation via sum-of-squares.
836 *CoRR*, abs/1711.11581, 2017. [1](#), [2](#), [4](#), [5](#)
- 837 [38] Pravesh K. Kothari and David Steurer. List-decodable mean estimation made simple. In
838 *Manuscript*, 2019. [2](#)
- 839 [39] Kevin A. Lai, Anup B. Rao, and Santosh Vempala. Agnostic estimation of mean and covariance.
840 In *FOCS*, pages 665–674. IEEE Computer Society, 2016. [1](#)
- 841 [40] Jean B. Lasserre. New positive semidefinite relaxations for nonconvex quadratic programs.
842 In *Advances in convex analysis and global optimization (Pythagorion, 2000)*, volume 54 of
843 *Nonconvex Optim. Appl.*, pages 319–331. Kluwer Acad. Publ., Dordrecht, 2001. [8](#)
- 844 [41] Monique Laurent. Sums of squares, moment matrices and optimization over polynomials. In
845 *Emerging applications of algebraic geometry*, pages 157–270. Springer, 2009. [10](#)
- 846 [42] Yuanzhi Li and Yingyu Liang. Learning mixtures of linear regressions with nearly optimal
847 complexity. In *Conference On Learning Theory, COLT 2018, Stockholm, Sweden, 6-9 July*
848 *2018.*, pages 1125–1144, 2018. [2](#), [4](#)
- 849 [43] Doron S Lubinsky. A Survey of Weighted Approximation for Exponential Weights. *arXiv*
850 *Mathematics e-prints*, page math/0701099, Jan 2007. [2](#), [18](#)
- 851 [44] Tengyu Ma, Jonathan Shi, and David Steurer. Polynomial-time tensor decompositions with
852 sum-of-squares. In *FOCS*, pages 438–446. IEEE Computer Society, 2016. [5](#), [8](#)
- 853 [45] RARD Maronna, R Douglas Martin, and Victor Yohai. *Robust statistics*. John Wiley & Sons,
854 Chichester. ISBN, 2006. [1](#)
- 855 [46] Ankur Moitra and Gregory Valiant. Settling the polynomial learnability of mixtures of gaussians.
856 In *FOCS*, pages 93–102. IEEE Computer Society, 2010. [4](#)
- 857 [47] Yurii Nesterov. Squared functional systems and optimization problems. In *High performance*
858 *optimization*, volume 33 of *Appl. Optim.*, pages 405–440. Kluwer Acad. Publ., Dordrecht, 2000.
859 [8](#)
- 860 [48] Pablo A Parrilo. *Structured semidefinite programs and semialgebraic geometry methods in*
861 *robustness and optimization*. PhD thesis, California Institute of Technology, 2000. [8](#)
- 862 [49] Adarsh Prasad, Arun Sai Suggala, Sivaraman Balakrishnan, and Pradeep Ravikumar. Robust
863 estimation via robust gradient estimation. *CoRR*, abs/1802.06485, 2018. [2](#)
- 864 [50] Mark Rudelson and Roman Vershynin. The Littlewood-Offord problem and invertibility of
865 random matrices. *Adv. Math.*, 218(2):600–633, 2008. [3](#)
- 866 [51] Hanie Sedghi, Majid Janzamin, and Anima Anandkumar. Provable tensor methods for learning
867 mixtures of generalized linear models. In *AISTATS*, volume 51 of *JMLR Workshop and*
868 *Conference Proceedings*, pages 1223–1231. JMLR.org, 2016. [2](#), [4](#)
- 869 [52] N. Z. Shor. Quadratic optimization problems. *Izv. Akad. Nauk SSSR Tekhn. Kibernet.*, (1):128–
870 139, 222, 1987. [8](#)

- 871 [53] Terence Tao and Van Vu. The Littlewood-Offord problem in high dimensions and a conjecture
872 of Frankl and Füredi. *Combinatorica*, 32(3):363–372, 2012. 3
- 873 [54] John W. Tukey. Mathematics and the picturing of data. pages 523–531, 1975. 1
- 874 [55] Eric W. Weisstein. Hermite number from mathworld. <http://mathworld.wolfram.com/HermiteNumber.html>. 20
875
- 876 [56] Wikipedia. Singleton bound. [https://en.wikipedia.org/wiki/Singleton_](https://en.wikipedia.org/wiki/Singleton_bound)
877 [bound](https://en.wikipedia.org/wiki/Singleton_bound). 17
- 878 [57] Xinyang Yi, Constantine Caramanis, and Sujay Sanghavi. Alternating Minimization for Mixed
879 Linear Regression. *arXiv e-prints*, page arXiv:1310.3745, Oct 2013. 2
- 880 [58] Kai Zhong, Prateek Jain, and Inderjit S. Dhillon. Mixed linear regression with multiple
881 components. In *NIPS*, pages 2190–2198, 2016. 2