
Statistical Query Lower Bounds for List-Decodable Linear Regression

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 We study the problem of list-decodable linear regression, where an adversary can
2 corrupt a majority of the examples. Specifically, we are given a set T of labeled
3 examples $(x, y) \in \mathbb{R}^d \times \mathbb{R}$ and a parameter $0 < \alpha < 1/2$ such that an α -fraction
4 of the points in T are i.i.d. samples from a linear regression model with Gaussian
5 covariates, and the remaining $(1 - \alpha)$ -fraction of the points are drawn from an
6 arbitrary noise distribution. The goal is to output a small list of hypothesis vectors
7 such that at least one of them is close to the target regression vector. Our main
8 result is a Statistical Query (SQ) lower bound of $d^{\text{poly}(1/\alpha)}$ for this problem. Our
9 SQ lower bound qualitatively matches the performance of previously developed
10 algorithms, providing evidence that current upper bounds for this task are nearly
11 best possible.

12 1 Introduction

13 1.1 Background and Motivation

14 Linear regression is one of the oldest and most fundamental statistical tasks with numerous applica-
15 tions in the sciences [RL87, Die01, McD09]. In the standard setup, the data are labeled examples
16 $(x^{(i)}, y^{(i)})$, where the examples (covariates) $x^{(i)}$ are i.i.d. samples from a distribution D_x on \mathbb{R}^d and
17 the labels $y^{(i)}$ are noisy evaluations of a linear function. More specifically, each label is of the form
18 $y^{(i)} = \beta \cdot x^{(i)} + \eta^{(i)}$, where $\eta^{(i)}$ is the observation noise, for an unknown target regression vector
19 $\beta \in \mathbb{R}^d$. The objective is to approximately recover the hidden regression vector. In this basic setting,
20 linear regression is well-understood. For example, under Gaussian distribution, the least-squares
21 estimator is known to be statistically and computationally efficient.

22 Unfortunately, classical efficient estimators inherently fail in the presence of even a very small
23 fraction of adversarially corrupted data. In several applications of modern data analysis, including
24 machine learning security [BNJT10, BNL12, SKL17, DKK⁺19] and exploratory data analysis, e.g.,
25 in biology [RPW⁺02, PLJD10, LAT⁺08], typical datasets contain arbitrary or adversarial outliers.
26 Hence, it is important to understand the algorithmic possibilities and fundamental limits of learning
27 and inference in such settings. Robust statistics focuses on designing estimators tolerant to a small
28 amount of contamination, where the outliers are the *minority* of the dataset. Classical work in this
29 field [HRRS86, HR09] developed robust estimators for various basic tasks, alas with exponential
30 runtime. More recently, a line of work in computer science, starting with [DKK⁺16, LRV16],
31 developed the first computationally efficient robust learning algorithms for various high-dimensional
32 tasks. Subsequently, there has been significant progress in algorithmic robust statistics by several
33 communities, see [DK19] for a survey on the topic.

34 In this paper, we study high-dimensional robust linear regression in the presence of a *majority* of
35 adversarial outliers. As we explain below, in several applications, asking for a minority of outliers

36 is too strong of an assumption. It is thus natural to ask what notion of learning can capture the
 37 regime when the clean data points (inliers) constitute the *minority* of the dataset. While outputting a
 38 *single* accurate hypothesis in this regime is information-theoretically impossible, one may be able
 39 to compute a *small list* of hypotheses with the guarantee that *at least one of them* is accurate. This
 40 relaxed notion is known as *list-decodable learning* [BBV08, CSV17], formally defined below.

41 **Definition 1.1.** (List-Decodable Learning) Given a parameter $0 < \alpha < 1/2$ and a distribution
 42 family \mathcal{D} on \mathbb{R}^d , the algorithm specifies $n \in \mathbb{Z}_+$ and observes n i.i.d. samples from a distribution
 43 $E = \alpha D + (1-\alpha)N$, where D is an unknown distribution in \mathcal{D} and N is arbitrary. We say D is
 44 the distribution of inliers, N is the distribution of outliers, and E is an $(1-\alpha)$ -corrupted version of
 45 D . Given sample access to an $(1-\alpha)$ -corrupted version of D , the goal is to output a “small” list of
 46 hypotheses \mathcal{L} at least one of which is (with high probability) close to the target parameter of D .

47 We note that a list of size $O(1/\alpha)$ typically suffices; an algorithm with a $\text{poly}(1/\alpha)$ sized list, or
 48 even a worse function of $1/\alpha$ (but independent of the dimension d) is also considered acceptable.

49 Natural applications of list-decodable learning include crowdsourcing, where a majority of partic-
 50 ipants could be unreliable [SVC16, MV18], and semi-random community detection in stochas-
 51 tic block models [CSV17]. List-decoding is also useful in the context of semi-verified learn-
 52 ing [CSV17, MV18], where a learner can audit a very small amount of trusted data. If the trusted
 53 dataset is too small to directly learn from, using a list-decodable learning procedure, one can pinpoint a
 54 candidate hypothesis consistent with the verified data. Importantly, list-decodable learning generalizes
 55 the task of learning mixture models, see, e.g., [DeV89, JY94, ZJD16, LL18, KC20, CLS20, DK20] for
 56 the case of linear regression studied here. Roughly speaking, by running a list-decodable estimation
 57 procedure with the parameter α equal to the smallest mixing weight, each true cluster of points is an
 58 equally valid ground-truth distribution, so the output list must contain candidate parameters close to
 59 each of the true parameters.

60 In list-decodable linear regression (the focus of this paper), D is a distribution on pairs (X, y) , where
 61 X is a standard Gaussian on \mathbb{R}^d , y is approximately a linear function of x , and the algorithm is asked
 62 to approximate the hidden regressor. The following definition specifies the distribution family \mathcal{D} of
 63 the inliers for the case of linear regression with Gaussian covariates.

64 **Definition 1.2.** (Gaussian Linear Regression) Fix $\sigma > 0$. For $\beta \in \mathbb{R}^d$, let D_β be the distribution
 65 over (X, y) , $X \in \mathbb{R}^d$, $y \in \mathbb{R}$, such that $X \sim \mathcal{N}(0, I_d)$ and $y = \beta^T X + \eta$, where $\eta \sim \mathcal{N}(0, \sigma^2)$
 66 independently of X . We define \mathcal{D} to be the set $\{D_\beta : \beta \in S'\}$ for some set $S' \subseteq \mathbb{R}^d$.

67 Recent algorithmic progress [KKK19, RY20] has been made on this problem using the SoS hierarchy.
 68 The guarantees in [KKK19, RY20] are very far from the information-theoretic limit in terms of
 69 sample complexity. In particular, they require $d^{\text{poly}(1/\alpha)}$ samples and time to obtain non-trivial error
 70 guarantees (see Table 1): [KKK19] obtains an error guarantee of $O(\sigma/\alpha)$ with a list of size $O(1/\alpha)$,
 71 whereas [RY20] obtains an error guarantee of $O(\sigma/\alpha^{3/2})$ with a list of size $(1/\alpha)^{O(\log(1/\alpha))}$.

72 On the other hand, as shown in this paper (see Theorem 1.4), $\text{poly}(d/\alpha)$ samples information-
 73 theoretically suffice to obtain near-optimal error guarantees. This raises the following natural
 74 question:

75 *What is the complexity of list-decodable linear regression?*
 76 *Are there efficient algorithms with significantly better sample-time tradeoffs?*

77 We study the above question in a natural and well-studied restricted model of computation, known as
 78 the Statistical Query (SQ) model [Kea98]. As the main result of this paper, we prove strong SQ lower
 79 bounds for this problem. Via a recently established equivalence [BBH⁺20], our SQ lower bound also
 80 implies low-degree testing lower bounds for this task. Our lower bounds can be viewed as evidence
 81 that current upper bounds for this problem may be qualitatively best possible.

82 Before we state our contributions in detail, we give some background on SQ algorithms. SQ
 83 algorithms are a broad class of algorithms that are only allowed to query expectations of bounded
 84 functions of the distribution rather than directly access samples. Formally, an SQ algorithm has
 85 access to the following oracle.

86 **Definition 1.3** (STAT Oracle). Let D be a distribution on \mathbb{R}^d . A statistical query is a bounded
 87 function $q : \mathbb{R}^d \rightarrow [-1, 1]$. For $\tau > 0$, the $\text{STAT}(\tau)$ oracle responds to the query q with a value v
 88 such that $|v - \mathbf{E}_{X \sim D}[q(X)]| \leq \tau$. We call τ the tolerance of the statistical query.

Table 1: The table summarizes the sample complexity, running time, and list size of the known list-decodable linear regression algorithms in order to obtain a $1/4$ -additive approximation to the hidden regression vector β in the setting of Theorem 1.5, i.e., when $\|\beta\|_2 \leq 1$ and σ is sufficiently small as a function of α : [KKK19] requires $\sigma = O(\alpha)$ and [RY20] requires $\sigma = O(\alpha^{3/2})$.

Algorithmic Result	Sample Size	Running Time	List size
Karmalkar-Klivans-Kothari [KKK19]	$(d/\alpha)^{O(1/\alpha^4)}$	$(d/\alpha)^{O(1/\alpha^8)}$	$O(1/\alpha)$
Raghavendra and Yau [RY20]	$d^{O(1/\alpha^4)}$	$d^{O(1/\alpha^8)}(1/\alpha)^{\log(1/\alpha)}$	$(1/\alpha)^{O(\log(1/\alpha))}$

89 The SQ model was introduced by Kearns [Kea98] in the context of supervised learning as a natural
90 restriction of the PAC model [Val84]. Subsequently, the SQ model has been extensively studied in a
91 plethora of contexts (see, e.g., [Fel16] and references therein). The class of SQ algorithms is rather
92 broad and captures a range of known supervised learning algorithms. More broadly, several known
93 algorithmic techniques in machine learning are known to be implementable using SQs. These include
94 spectral techniques, moment and tensor methods, local search (e.g., Expectation Maximization), and
95 many others (see, e.g., [FGR⁺17, FGV17]).

96 1.2 Our Results

97 We start by showing that $\text{poly}(d/\alpha)$ samples are sufficient to obtain a near-optimal error estimator,
98 albeit with a computationally inefficient algorithm.

99 **Theorem 1.4** (Information-Theoretic Bound). *There is a (computationally inefficient) list-decoding*
100 *algorithm for Gaussian linear regression that uses $O(d/\alpha^3)$ samples, returns a list of $O(1/\alpha)$ many*
101 *hypothesis vectors, and has ℓ_2 -error guarantee of $O((\sigma/\alpha)\sqrt{\log(1/\alpha)})$. Moreover, if the dimension*
102 *d is sufficiently large, any list-decoding algorithm that outputs a list of size $\text{poly}(1/\alpha)$ must have*
103 *ℓ_2 -error at least $\Omega((\sigma/\alpha)/\sqrt{\log(1/\alpha)})$.*

104 Due to space limitations, the proof of Theorem 1.4 is deferred to the supplementary material (see
105 Theorems D.2 and D.4).

106 Our main result is a strong SQ lower bound for the list-decodable Gaussian linear regression problem.
107 We establish the following theorem (see Theorem 2.1 for a more detailed formal statement).

108 **Theorem 1.5** (SQ Lower Bound). *Assume that the dimension $d \in \mathbb{Z}_+$ is sufficiently large and*
109 *consider the problem of list-decodable linear regression, where the fraction of inliers is $\alpha \in (0, 1/2)$,*
110 *the regression vector $\beta \in \mathbb{R}^d$ has norm $\|\beta\|_2 \leq 1$, and the additive noise has standard deviation*
111 *$\sigma \leq \alpha$. Then any SQ algorithm that returns a list \mathcal{L} of candidate vectors containing a $\hat{\beta}$ such that*
112 *$\|\hat{\beta} - \beta\|_2 \leq 1/4$ does one of the following: (i) it uses at least one query with tolerance at most*
113 *$d^{-\Omega(1/\sqrt{\alpha})}/\sigma$, (ii) it makes $2^{d^{\Omega(1)}}$ queries, or (iii) it returns a list of size $|\mathcal{L}| = 2^{d^{\Omega(1)}}$.*

114 Informally speaking, Theorem 1.5 shows that no SQ algorithm can approximate β to constant
115 accuracy with a sub-exponential in $d^{\Omega(1)}$ size list and sub-exponential in $d^{\Omega(1)}$ many queries, unless
116 using queries of very small tolerance – that would require at least $\sigma d^{\Omega(1/\sqrt{\alpha})}$ samples to simulate.
117 For σ not too small, e.g., $\sigma = \text{poly}(\alpha)$, in view of Theorem 1.4, this result can be viewed as an
118 information-computation tradeoff for the problem, within the class of SQ algorithms.

119 A conceptual implication of Theorem 1.5 is that list-decodable linear regression is harder (within
120 the class of SQ algorithms) than the related problem of learning mixtures of linear regressions
121 (MLR). Recent work [DK20] gave an algorithm (easily implementable in SQ) for learning MLR
122 with k equal weight separated components (under Gaussian covariates) with sample complexity and
123 running time $k^{\text{poly}(\log(k))}$, i.e., *quasi-polynomial* in k . Recalling that one can reduce k -MLR (with
124 well-separated components) to list-decodable linear regression for $\alpha = 1/k$, Theorem 1.5 implies
125 that the aforementioned algorithmic result cannot be obtained via such a reduction.

126 **Remark 1.6.** While the main focus of this work is on the SQ model, our result has immediate
127 implications to a related popular restricted computational model — that of low-degree (polynomial)
128 algorithms [HS17, HKP⁺17, Hop18]. Recent work [BBH⁺20] established that (under certain as-

129 sumptions) an SQ lower bound also implies a qualitatively similar lower bound in the low-degree
 130 model. We leverage this connection to show a similar lower bound in this model (see Appendix F).

131 1.3 Overview of Techniques

132 In this section, we provide a detailed overview of our SQ lower bound construction. We recall that
 133 there exists a general methodology for establishing SQ lower bounds via an appropriate complexity
 134 measure, known as SQ dimension. Several related notions of SQ dimension exist in the literature,
 135 see, e.g., [BFJ⁺94, FGR⁺17, Fel17]. Here we focus on the framework introduced in [FGR⁺17]
 136 for search problems over distributions, which is more natural in our setting. A lower bound on the
 137 SQ dimension of a search problem provides an unconditional lower bound on the SQ complexity
 138 of the problem. Roughly speaking, for a notion of correlation between distributions in our family
 139 \mathcal{D} (Definition 1.8), establishing an SQ lower bound amounts to constructing a large cardinality
 140 sub-family $\mathcal{D}' \subseteq \mathcal{D}$ such that every pair of distributions in \mathcal{D}' are nearly uncorrelated with respect to
 141 a given reference distribution R (see Definition 1.10 and Lemma 1.11).

142 A general framework for constructing SQ-hard families of distributions was introduced in [DKS17],
 143 which showed the following: Let the reference distribution R be $\mathcal{N}(0, I)$ and A be a univariate
 144 distribution whose low-degree moments match those of the standard Gaussian (and which satisfies
 145 an additional mild technical condition). Let $P_{A,v}$ be the distribution that is a copy of A in the v -
 146 direction and standard Gaussian in the orthogonal complement (Definition 1.12). Then the distribution
 147 family $\{P_{A,v}\}_{v \in S}$, where S is a set of nearly orthogonal unit vectors, satisfies the pairwise nearly
 148 uncorrelated property (Lemma 1.13), and is therefore SQ-hard to learn.

149 Unfortunately, the [DKS17] framework does not suffice in the supervised setting of the current paper
 150 for the following reason: The joint distribution over labeled examples (X, y) in our setting does not
 151 possess the symmetry properties required for moment-matching with the reference $R = \mathcal{N}(0, I)$ to
 152 be possible. Specifically, the behavior of y will necessarily be somewhat different than the behavior of
 153 X . To circumvent this issue, we leverage an idea from [DKS19]. The high-level idea is to construct
 154 distributions E_v on (X, y) such that for any fixed value y_0 of y , the conditional distribution of
 155 $X \mid y = y_0$ under E_v is of the form $P_{A,v}$ described above, where A is replaced with some A_{y_0} .

156 We further explain this modified construction. Note that E_v should be of the form $\alpha D_v + (1-\alpha)N_v$,
 157 where D_v is the inlier distribution (corresponding to the clean samples from the linear regression
 158 model) and N_v is the outlier (noise) distribution. To understand what properties our distribution
 159 should satisfy, we start by looking at the inlier distribution D . By definition, for $(X, y) \sim D$, we
 160 have that $y = \beta^T X + \eta$, where $X \sim \mathcal{N}(0, I)$ and $\eta \sim \mathcal{N}(0, \sigma^2)$ is independent of X . A good
 161 place to start here is to understand the distribution of X conditioned on $y = y_0$, for some y_0 , under
 162 D . It is not hard to show (Fact 2.3) that this conditional distribution is already of the desired form
 163 $P_{A,\beta}$: it is a product of a $(d-1)$ -dimensional standard Gaussian in directions orthogonal to β ,
 164 while in the β -direction it is a much narrower Gaussian with mean proportional to y_0 . To establish
 165 our SQ-hardness result, we would like to mix this conditional distribution with a carefully selected
 166 outlier distribution $N \mid y = y_0$, such that the resulting mixture $E \mid y = y_0$ matches many of its
 167 low-degree moments with the standard Gaussian in the β -direction, while being standard Gaussian
 168 in the orthogonal directions. In the setting of minority of outliers, [DKS19] was able to provide an
 169 explicit formula for N and match *three* moments to show an SQ lower bound of $\Omega(d^2)$. The main
 170 technical difficulty in our paper is that, in order to prove the desired SQ lower bound of $\Omega(d^{\text{poly}(1/\alpha)})$,
 171 we need to match $\text{poly}(1/\alpha)$ many moments. We explain how to achieve this below.

172 Here we take a different approach and establish the existence of the desired outlier distribution
 173 $N \mid y = y_0$ in a non-constructive manner. We note that our problem is an instance of the moment-
 174 matching problem, where given a sequence of real numbers, the goal is to decide whether a distribution
 175 exists having that sequence as its low-degree moments. At a high-level, we leverage classical results
 176 that tackle this general question by formulating a linear program (LP) and using LP-duality to derive
 177 necessary and sufficient feasibility conditions (see [KS53] and Theorem 3.1). This moment-matching
 178 via LP duality approach is fairly general, but stumbles upon two technical obstacles in our setting.

179 The first technical issue is that our final distributions E_v on (X, y) need to have bounded χ^2 -
 180 divergence with respect to the reference distribution, since the pairwise correlations scale with this
 181 quantity (see Lemma 1.13). To guarantee this, we can ensure that the outlier distribution in the
 182 β -direction is in fact equal to the convolution of a distribution with bounded support with a narrow

183 Gaussian: (i) The contraction property of this convolution operator means that it can only reduce the
 184 χ^2 -divergence, and (ii) the bounded support can be used in combination with tail-bounds on Hermite
 185 polynomials (Lemma 3.6) to bound from above the contribution to the χ^2 -divergence of each Hermite
 186 coefficient of our distribution (Lemma 2.6). These additional constraints necessitate a modification to
 187 the moment-matching problem, but it can still be readily analyzed (Theorem 2.5).

188 The second and more complicated issue involves the fraction of outliers, i.e., the parameter “ $1-\alpha$ ”.
 189 Unfortunately, it is easy to see that the fraction of outliers necessary to make the conditional
 190 distributions match the desired number of moments must necessarily go to 1 as $|y|$ goes to infinity:
 191 As $|y|$ gets bigger, the conditional distribution of inliers moves further away from $\mathcal{N}(0, I)$ (Fact 2.3)
 192 and thus needs to be mixed more heavily with outliers to be corrected. This is a significant problem,
 193 since by definition we can only afford to use a $(1-\alpha)$ -fraction of outliers overall. To handle this issue,
 194 we consider a reference distribution R on (X, y) that has much heavier tails in y than the distribution
 195 of inliers has. This essentially means that as $|y|$ gets large, the conditional probability that a sample
 196 is an outlier gets larger and larger. This is balanced by having slightly lower fraction of outliers for
 197 smaller values of $|y|$, in order to ensure that the total fraction of outliers is still at most $1-\alpha$. To
 198 address this issue, we leverage the fact that the probability that a clean sample has large value of $|y|$
 199 is very small. Consequently, we can afford to make the error rates for such y quite large without
 200 increasing the overall probability of error by very much.

201 1.4 Preliminaries

202 **Notation** We use \mathbb{N} to denote natural numbers and \mathbb{Z}_+ to denote positive integers. For $n \in \mathbb{Z}_+$ we
 203 denote $[n] \stackrel{\text{def}}{=} \{1, \dots, n\}$ and use S^{d-1} for the d -dimensional unit sphere. We denote by $\mathbf{1}(\mathcal{E})$ the
 204 indicator function of the event \mathcal{E} . We use I_d to denote the $d \times d$ identity matrix. For a random
 205 variable X , we use $\mathbf{E}[X]$ for its expectation. For $m \in \mathbb{Z}_+$, the m -th moment of X is defined as
 206 $\mathbf{E}[X^m]$. We use $\mathcal{N}(\mu, \Sigma)$ to denote the Gaussian distribution with mean μ and covariance matrix Σ .
 207 We let ϕ denote the pdf of the one-dimensional standard Gaussian. When D is a distribution, we use
 208 $X \sim D$ to denote that the random variable X is distributed according to D . For a vector $x \in \mathbb{R}^d$, we
 209 let $\|x\|_2$ denote its ℓ_2 -norm. For $y \in \mathbb{R}$, we denote by δ_y the Dirac delta distribution at y , i.e., the
 210 distribution that assigns probability mass 1 to the single point $y \in \mathbb{R}$ and zero elsewhere. When there
 211 is no confusion, we will use the same letters for distributions and their probability density functions.

212 **Ornstein-Uhlenbeck Operator** For a $\rho > 0$, we define the *Gaussian noise* (or *Ornstein-Uhlenbeck*)
 213 operator U_ρ as the operator that maps a distribution F on \mathbb{R} to the distribution of the random variable
 214 $\rho X + \sqrt{1-\rho^2}Z$, where $X \sim F$ and $Z \sim \mathcal{N}(0, 1)$ independently of X .

215 **Background on the SQ Model** We provide the basic definitions and facts that we use.

216 **Definition 1.7** (Search problems over distributions). Let \mathcal{D} be a set of distributions over \mathbb{R}^d , \mathcal{F} be
 217 a set called solutions, and $\mathcal{Z} : \mathcal{D} \rightarrow 2^{\mathcal{F}}$ be a map that assigns sets of solutions to distributions of
 218 \mathcal{D} . The *distributional search problem* \mathcal{Z} over \mathcal{D} and \mathcal{F} is to find a valid solution $f \in \mathcal{Z}(D)$ given
 219 statistical query oracle access to an unknown $D \in \mathcal{D}$.

220 The hardness of these problems is conveniently captured by the SQ dimension. For this, we first need
 221 to define the notion of correlation between distributions.

222 **Definition 1.8** (Pairwise Correlation). The pairwise correlation of two distributions with probability
 223 density functions $D_1, D_2 : \mathbb{R}^d \rightarrow \mathbb{R}_+$ with respect to a reference distribution with density $R : \mathbb{R}^d \rightarrow$
 224 \mathbb{R}_+ , where the support of R contains the supports of D_1 and D_2 , is defined as $\chi_R(D_1, D_2) \stackrel{\text{def}}{=} \int_{\mathbb{R}^d} D_1(x)D_2(x)/R(x) dx - 1$. When $D_1 = D_2$, the pairwise correlation becomes the same as the
 225 χ^2 -divergence between D_1 and R , i.e., $\chi^2(D_1, R) \stackrel{\text{def}}{=} \int_{\mathbb{R}^d} D_1^2(x)/R(x) dx - 1$.

227 **Definition 1.9.** For $\gamma, \beta > 0$, the set of distributions $\mathcal{D} = \{D_1, \dots, D_m\}$ is called (γ, β) -correlated
 228 relative to the distribution R if $|\chi_R(D_i, D_j)| \leq \gamma$, if $i \neq j$, and $|\chi_R(D_i, D_j)| \leq \beta$ otherwise.

229 The statistical dimension of a search problem is based on the largest set of (γ, β) -correlated distribu-
 230 tions assigned to each solution.

231 **Definition 1.10** (Statistical Dimension). For $\gamma, \beta > 0$, a search problem \mathcal{Z} over a set of solutions
 232 \mathcal{F} and a class \mathcal{D} of distributions over X , we define the *statistical dimension* of \mathcal{Z} , denoted by
 233 $\text{SD}(\mathcal{Z}, \gamma, \beta)$, to be the largest integer m such that there exists a reference distribution R over X and

234 a finite set of distributions $\mathcal{D}_R \subseteq \mathcal{D}$ such that for any solution $f \in \mathcal{F}$, the set $\mathcal{D}_f = \mathcal{D}_R \setminus \mathcal{Z}^{-1}(f)$ is
 235 (γ, β) -correlated relative to R and $|\mathcal{D}_f| \geq m$.

236 **Lemma 1.11** (Corollary 3.12 in [FGR⁺17]). *Let \mathcal{Z} be a search problem over a set of solutions \mathcal{F} and
 237 a class of distributions \mathcal{D} over \mathbb{R}^d . For $\gamma, \beta > 0$, let $s = \text{SD}(\mathcal{Z}, \gamma, \beta)$ be the statistical dimension
 238 of the problem. For any $\gamma' > 0$, any SQ algorithm for \mathcal{Z} requires either $s\gamma'/(\beta - \gamma)$ queries or at
 239 least one query to $\text{STAT}(\sqrt{\gamma + \gamma'})$ oracle.*

240 We continue by recalling the machinery from [DKS17] that will be used for our construction.

241 **Definition 1.12** (High-Dimensional Hidden Direction Distribution). For a unit vector $v \in \mathbb{R}^d$ and a
 242 distribution A on the real line with probability density function $A(x)$, define $P_{A,v}$ to be a distribution
 243 over \mathbb{R}^d , where $P_{A,v}$ is the product distribution whose orthogonal projection onto the direction of v is
 244 A , and onto the subspace perpendicular to v is the standard $(d-1)$ -dimensional normal distribution.
 245 That is, $P_{A,v}(x) := A(v^T x) \phi_{\perp v}(x)$, where $\phi_{\perp v}(x) = \exp(-\|x - (v^T x)v\|_2^2/2) / (2\pi)^{(d-1)/2}$.

246 The distributions $\{P_{A,v}\}$ defined above are shown to be nearly uncorrelated as long as the directions
 247 where A is embedded are pairwise nearly orthogonal.

248 **Lemma 1.13** (Lemma 3.4 in [DKS17]). *Let $m \in \mathbb{Z}_+$. Let A be a distribution over \mathbb{R} that agrees
 249 with the first m moments of $\mathcal{N}(0, 1)$. For any v , let $P_{A,v}$ denote the distribution from Definition 1.12.
 250 For all $v, u \in \mathbb{R}^d$, we have that $\chi_{\mathcal{N}(0, I_d)}(P_{A,v}, P_{A,u}) \leq |u^T v|^{m+1} \chi^2(A, \mathcal{N}(0, 1))$.*

251 The following result shows that there are exponentially many nearly-orthogonal unit vectors.

252 **Lemma 1.14** (see, e.g., Lemma 3.7 of [DKS17]). *For any $0 < c < 1/2$, there is a set S , of at least
 253 $2^{\Omega(d^c)}$ unit vectors in \mathbb{R}^d , such that for each pair of distinct $v, v' \in S$, it holds $|v^T v'| \leq O(d^{c-1/2})$.*

254 2 Main Result: Proof of Theorem 1.5

255 In this section, we present the main result of this paper: SQ hardness of list-decodable linear regression
 256 (Definitions 1.1 and 1.2). We consider the setting when β has norm less than 1, i.e., $\beta = \rho v$ for
 257 $v \in \mathcal{S}^{d-1}$ and $\rho \in (0, 1)$.¹ Note that the marginal distribution of the labels is $\mathcal{N}(0, \sigma_y^2)$, where
 258 $\sigma_y^2 = \rho^2 + \sigma^2$. We ensure that the labels y have unit variance by using $\sigma^2 = 1 - \rho^2$. Specifically, the
 259 choice of parameters will be such that obtaining a $\rho/2$ -additive approximation of the regressor β is
 260 possible information-theoretically with $\text{poly}(d/\alpha)$ samples (cf. Appendix D.1), but the complexity of
 261 any SQ algorithm for the task must necessarily be at least $d^{\text{poly}(1/\alpha)}/\sigma$. We show the following more
 262 detailed statement of Theorem 1.5.

263 **Theorem 2.1** (SQ Lower Bound). *Let $c \in (0, 1/2)$, $d \in \mathbb{Z}_+$ with $d = 2^{\Omega(1/(1/2-c))}$, $\alpha \in (0, 1/2)$,
 264 $\rho \in (0, 1)$, $\sigma^2 = 1 - \rho^2$, and $m \in \mathbb{Z}_+$ with $m \leq c_1/\sqrt{\alpha}$ for some sufficiently small constant
 265 $c_1 > 0$. Any list-decoding algorithm that, given statistical query access to a $(1-\alpha)$ -corrupted
 266 version of the distribution described by the model of Definition 1.2 with $\beta = \rho v$ for $v \in \mathcal{S}^{d-1}$,
 267 returns a list \mathcal{L} of hypotheses vectors that contains a $\hat{\beta}$ such that $\|\hat{\beta} - \beta\|_2 \leq \rho/2$, does one of the
 268 following: (i) it uses at least one query to $\text{STAT}\left(\Omega(d)^{-(2m+1)(1/4-c/2)} e^{O(m)} / \sqrt{1 - \rho^2}\right)$, (ii) it
 269 makes $2^{\Omega(d^c)} d^{-(2m+1)(1/2-c)}$ many queries, or (iii) it returns a list \mathcal{L} of size $2^{\Omega(d^c)}$.*

270 In the rest of this section, we will explain the hard-to-learn construction for our SQ lower bound, i.e.,
 271 a set of distributions with large statistical dimension. The proof would then follow from Lemma 1.11.
 272 We begin by describing additional notation that we will use.

273 **Notation:** As $\beta = \rho v$ for a fixed ρ , we will slightly abuse notation by using $D_v(x, y)$ to denote
 274 the joint distribution of the inliers and we use $E_v(x, y)$ to denote the $(1-\alpha)$ -corrupted version of
 275 $D_v(x, y)$. To avoid using multiple subscripts, we use $D_v(x|y)$ to denote the conditional distribution
 276 of $X|y$ according to the distribution D_v and similarly for the other distributions. In addition, we use
 277 $D_v(y)$ to denote the marginal distribution of y under D_v and similarly for other distributions.

278 Following the general construction of [DKS17], we will specify a *reference* joint distribution $R(x, y)$
 279 where X and y are independent, and $X \sim \mathcal{N}(0, I_d)$. We will find a marginal distribution $R(y)$ and a
 280 way to add the outliers so that the following hold for each E_v (where $m = \Theta(1/\sqrt{\alpha})$):

¹This is a standard assumption and considered by existing works [KKK19, RY20] (cf. Remark B.5).

- 281 (I) E_v is indeed a valid distribution of (X, y) in our corruption model (i.e., can be written as a
 282 mixture $\alpha D_v(x, y) + (1-\alpha)N_v(x, y)$ for some noise distribution N_v). Moreover, the marginal
 283 of E_v on the labels, $E_v(y)$, coincides with $R(y)$.
- 284 (II) For every $y \in \mathbb{R}$, the conditional distribution $E_v(x|y)$ is of the form $P_{A_y, v}$ of Definition 1.12,
 285 with A_y being a distribution that matches the first $2m$ moments with $\mathcal{N}(0, 1)$.²
- 286 (III) For A_y defined above, $\mathbf{E}_{y \sim R(y)}[\chi^2(A_y, \mathcal{N}(0, 1))]$ is bounded.

287 We first briefly explain why a construction satisfying the above properties suffices to prove our main
 288 theorem (postponing a formal proof for the end of this section). We start by noting the following
 289 decomposition (proved in Appendix B).

290 **Lemma 2.2.** *For $u, v \in \mathcal{S}^{d-1}$, if E_u and E_v have the same marginals $R(y)$ on the labels, they satisfy*
 291 $\chi_{R(x, y)}(E_v(x, y), E_u(x, y)) = \mathbf{E}_{y \sim R(y)}[\chi_{\mathcal{N}(0, I_d)}(E_v(x|y), E_u(x|y))]$.

292 Using the decomposition in Lemma 2.2 for E_u and E_v satisfying Property (II), Lemma 1.13 implies
 293 that $|\chi_{R(x, y)}(E_v(x, y), E_u(x, y))| \leq |u^T v|^{2m+1} \mathbf{E}_{y \sim R(y)}[\chi^2(A_y, \mathcal{N}(0, 1))]$. Letting $\mathcal{D} = \{E_v :$
 294 $v \in S\}$, where S is the set of nearly uncorrelated unit vectors from Lemma 1.14, we get that
 295 \mathcal{D} is (γ, b) -correlated relative to R , for $b = \mathbf{E}_{y \sim R(y)}[\chi^2(A_y, \mathcal{N}(0, 1))]$ and $\gamma \leq d^{-\Omega(m)b}$. As
 296 $|S| = 2^{\Omega(d^c)}$, b is bounded, and the list size is much smaller than $|S|$, we can show that the statistical
 297 dimension of the list-decodable linear regression is large.

298 Thus, in the rest of the section we focus on showing that such a construction exists. We first note
 299 that according to our linear model of Definition 1.2, the conditional distribution of X given y for the
 300 inliers is Gaussian with unit variance in all but one direction (see Appendix B for a proof).

301 **Fact 2.3.** *Fix $\rho > 0$, $v \in \mathcal{S}^{d-1}$, and consider the regression model of Definition 1.2 with $\beta = \rho v$.*
 302 *Then the conditional distribution $X|y$ of the inliers is $\mathcal{N}(y\rho v, I_d - \rho^2 v v^T)$, i.e., independent standard*
 303 *Gaussian in all directions perpendicular to v and $\mathcal{N}(\rho y, 1 - \rho^2)$ in the direction of v .*

304 Since Fact 2.3 states that $D_v(x|y)$ is already of the desired form (standard normal in all directions
 305 perpendicular to v and $\mathcal{N}(y\rho, 1 - \rho^2)$ in the direction of v), the problem becomes one-dimensional.
 306 More specifically, for every $y \in \mathbb{R}$, we need to find a one-dimensional distribution Q_y and appropriate
 307 values $\alpha_y \in [0, 1]$ such that the mixture $A_y = \alpha_y \mathcal{N}(y\rho, 1 - \rho^2) + (1-\alpha_y)Q_y$ matches the first
 308 $2m$ moments with $\mathcal{N}(0, 1)$. Then, multiplying by $\phi_{\perp v}$ (which denotes the contribution of the
 309 space orthogonal to v to the density of standard Gaussian, as defined in Definition 1.12) yields
 310 the d -dimensional mixture distribution $\alpha_y D_v(x|y) + (1-\alpha_y)Q_y(v^T x) \phi_{\perp v}(x)$. We show that an
 311 appropriate selection of α_y can ensure that this is a valid distribution for our contamination model.

312 **Lemma 2.4.** *Let R be a distribution on pairs $(x, y) \in \mathbb{R}^{d+1}$ such that $\alpha_y := \alpha D_v(y)/R(y) \in [0, 1]$*
 313 *for all $y \in \mathbb{R}$. Suppose that for every $y \in \mathbb{R}$ there exists a univariate distribution Q_y such that*
 314 *$A_y := \alpha_y \mathcal{N}(y\rho, 1 - \rho^2) + (1-\alpha_y)Q_y$ matches the first $2m$ moments with $\mathcal{N}(0, 1)$. If the distribution*
 315 *of the outliers is $N_v(x, y) = ((1-\alpha_y)/(1-\alpha))Q_y(v^T x) \phi_{\perp v}(x)R(y)$, Properties (I) and (II) hold.*

316 The proof of Lemma 2.4 is included in Appendix B. We will choose the reference distribution $R(x, y)$
 317 to have $X \sim \mathcal{N}(0, I_d)$ and $y \sim \mathcal{N}(0, 1/\alpha)$ independently, which makes the corresponding value of
 318 α_y to be $\alpha_y = \alpha D_v(y)/R(y) = \sqrt{\alpha} \exp(-y^2(1-\alpha)/2)$. This satisfies the condition in Lemma 2.4
 319 that $\alpha_y \in [0, 1]$. Our choice of $R(y) \sim \mathcal{N}(0, 1/\alpha)$ is informed by Properties (II) and (III), and will be
 320 used later on in the proofs of Theorem 2.5 and Lemma 2.6 (also see the last paragraph of Section 1.3
 321 for more intuition). Going back to our goal, i.e., making $A_y = \alpha_y \mathcal{N}(y\rho, 1 - \rho^2) + (1-\alpha_y)Q_y$ match
 322 moments with $\mathcal{N}(0, 1)$, we will argue that it suffices to only look for Q_y of the specific form $U_\rho F_y$,
 323 where U_ρ is the Ornstein-Uhlenbeck operator. This suffices because $U_\rho \delta_y = \mathcal{N}(y\rho, 1 - \rho^2)$ and the
 324 operator U_ρ preserves the moments of a distribution if they match with $\mathcal{N}(0, 1)$ (see Lemma 2.6 (i)
 325 below). Letting $A_y = U_\rho(\alpha_y \delta_y + (1 - \alpha_y)F_y)$, the new goal is to show that the argument of U_ρ
 326 matches moments with $\mathcal{N}(0, 1)$. We show the following structural result:

327 **Theorem 2.5.** *Let $y \in \mathbb{R}$, $B \in \mathbb{R}$, $\alpha \in (0, 1/2)$, and define $\alpha_y := \sqrt{\alpha} \exp(-y^2(1-\alpha)/2)$. For any*
 328 *$m \in \mathbb{Z}_+$ such that $m \leq C_1/\sqrt{\alpha}$ and $B \geq C_2\sqrt{m}$, with $C_1 > 0$ being a sufficiently small constant*
 329 *and C_2 being a sufficiently large constant, there exists a distribution F_y that satisfies the following:*

- 330 1. *The mixture distribution $\alpha_y \delta_y + (1 - \alpha_y)F_y$ matches the first $2m$ moments with $\mathcal{N}(0, 1)$.*

²We use even number of moments for simplicity. The analysis would slightly differ for odd number.

331 2. F_y is a discrete distribution supported on at most $2m + 1$ points, all of which lie in $[-B, B]$.

332 The proof of Theorem 2.5 is the bulk of the technical work of this paper and is deferred to Section 3. As
 333 mentioned before, applying U_ρ preserves the required moment-matching property. More crucially, it
 334 allows us to bound the χ^2 -divergence: the following result bounds $\chi^2(A_y, \mathcal{N}(0, 1))$ using contraction
 335 properties of U_ρ , tail bounds of Hermite polynomials, and the discreteness of F_y .

336 **Lemma 2.6.** *In the setting of Theorem 2.5, let $\rho > 0$ and $Q_y = U_\rho F_y$. Then the following holds*
 337 *for the mixture $A_y = \alpha_y \mathcal{N}(y\rho, 1 - \rho^2) + (1 - \alpha_y)Q_y$: (i) A_y matches the first $2m$ moments with*
 338 *$\mathcal{N}(0, 1)$, and (ii) $\chi^2(A_y, \mathcal{N}(0, 1)) \leq \alpha O(e^{y^2(\alpha-1/2)})/(1 - \rho^2) + O(e^{B^2/2})/(1 - \rho^2)$.*

339 We prove Lemma 2.6 in Appendix B. We are now ready to sketch the proof of Theorem 2.1 (see
 340 Appendix B for the detailed proof).

341 *Proof Sketch of Theorem 2.1.* Consider the search problem \mathcal{Z} , where \mathcal{D} is the set of all distributions
 342 E_v satisfying properties (I),(II), and (III) (let $\beta(v) = \rho v$ be the corresponding regressors). For each
 343 E_v , the corresponding solution set is defined to consist of all lists \mathcal{L} of size ℓ having one element
 344 that is $(\rho/2)$ -close to $\beta(v)$. Let the subset $\mathcal{D}_R = \{E_v\}_{v \in S}$, for S being the set of nearly orthogonal
 345 vectors of Lemma 1.14. Since $|u^T v| \leq O(d^{c-1/2})$ for any distinct $u, v \in S$ and $d = 2^{\Omega(1/(1/2-c))}$,
 346 for any vector w , at most one element of S can be $(\rho/2)$ -close to w . Thus, for any list \mathcal{L} of size
 347 $\ell = |S|/2$, $|\mathcal{D}_R \setminus \mathcal{Z}^{-1}(\mathcal{L})| \geq |S| - \ell \geq 2^{\Omega(d^c)}$. Using Lemmas 2.2 and 1.13 along with the χ^2 -bound
 348 of Lemma 2.6, we get that \mathcal{D}_R is (γ, b) -correlated with respect to R , for $b := e^{O(m)}/(1 - \rho^2)$ and
 349 $\gamma := \Omega(d)^{-(2m+1)(1/2-c)}$. An application of Lemma 1.11 completes the proof. \square

340 3 Duality for Moment Matching: Proof of Theorem 2.5

351 We now prove the existence of a bounded distribution F_y such that the mixture $\alpha_y \delta_y + (1 - \alpha_y)F_y$
 352 matches the first $2m$ moments with $\mathcal{N}(0, 1)$. The proof follows a non-constructive argument based
 353 on the duality between the space of moments and the space of non-negative polynomials.

354 Let $B > 0$ and $m \in \mathbb{Z}_+$. Let $\mathcal{P}(m)$ denote the class of all polynomials $p : \mathbb{R} \rightarrow \mathbb{R}$ with
 355 degree at most m . Let $\mathcal{P}^{\geq 0}(2m, B)$ be the class of polynomials that can be represented in either
 356 the form $p(t) = (\sum_{i=0}^m a_i t^i)^2$ or the form $p(t) = (B^2 - t^2)(\sum_{i=0}^{m-1} b_i t^i)^2$. The intuition for
 357 $\mathcal{P}^{\geq 0}(2m, B)$ is that every polynomial of degree at most $2m$ that is non-negative in $[-B, B]$ can
 358 be written as a finite sum of polynomials from $\mathcal{P}^{\geq 0}(2m, B)$. By slightly abusing notation, for
 359 a polynomial $p(t) = \sum_{i=0}^m p_i t^i$, we also use p to denote the vector in \mathbb{R}^{m+1} consisting of the
 360 coefficients (p_0, \dots, p_m) . The following classical result characterizes when a vector is realizable as
 361 the moment sequence of a distribution with support in $[-B, B]$ (for simplicity, we focus on matching
 362 an even number of moments in the rest of this section).

363 **Theorem 3.1** (Theorem 16.1 of [KS53]). *Let $B > 0$, $k \in \mathbb{Z}_+$, and $x = (x_0, x_1, \dots, x_{2k}) \in \mathbb{R}^{2k+1}$*
 364 *with $x_0 = 1$. There exists a distribution with support in $[-B, B]$ having as its first $2k$ moments the*
 365 *sequence (x_1, \dots, x_{2k}) if and only if for all $p \in \mathcal{P}^{\geq 0}(2k, B)$ it holds that $\sum_{i=0}^{2k} x_i p_i \geq 0$.*

366 As we require the distribution to be discrete, we prove the following result using Theorem 3.1:

367 **Proposition 3.2.** *Fix $y \in \mathbb{R}$, $\alpha_y \in (0, 1)$, $B > 0$, and $m \in \mathbb{Z}_+$. There exists a discrete distribution*
 368 *F_y supported on at most $2m + 1$ points in $[-B, B]$ such that $\alpha_y \delta_y + (1 - \alpha_y)F_y$ matches the first*
 369 *$2m$ moments with $\mathcal{N}(0, 1)$ if and only if $\mathbf{E}_{X \sim \mathcal{N}(0,1)}[p(X)] \geq \alpha_y p(y)$ for all $p \in \mathcal{P}^{\geq 0}(2m, B)$.*

370 The proof of Proposition 3.2 is deferred to Appendix C.1. To prove Theorem 2.5, we need to establish
 371 the condition of Proposition 3.2. To this end, we first need the following two technical lemmas, whose
 372 proofs are sketched towards the end of this section (for detailed proofs see Sections C.2 and C.3).

373 **Lemma 3.3.** *Let $m \in \mathbb{Z}_+$. If $B \geq C\sqrt{m}$ for some sufficiently large constant $C > 0$, then for every*
 374 *$q \in \mathcal{P}(m)$, it holds that $B^2 \mathbf{E}_{X \sim \mathcal{N}(0,1)}[q^2(X)] \geq 2 \mathbf{E}_{X \sim \mathcal{N}(0,1)}[X^2 q^2(X)]$.*

375 **Lemma 3.4.** *Let $y \in \mathbb{R}$, $\alpha \in (0, 1/2)$, $m \in \mathbb{Z}_+$, and $\alpha_y = \sqrt{\alpha} \exp(-y^2(1 - \alpha)/2)$. Sup-*
 376 *pose $m \leq C/\sqrt{\alpha}$ for some sufficiently small constant $C > 0$. Then for all $r \in \mathcal{P}(m)$, $r \neq 0$:*
 377 *$r^2(y)/(\mathbf{E}_{X \sim \mathcal{N}(0,1)}[r^2(X)]) \leq 1/(2\alpha_y)$.*

378 *Proof of Theorem 2.5.* By Proposition 3.2, it remains to show that if $B \geq C_2\sqrt{m}$, then the condition
 379 $\mathbf{E}_{X \sim \mathcal{N}(0,1)}[p(X)] \geq \alpha_y p(y)$ holds for all $p \in \mathcal{P}^{\geq 0}(2m, B)$. Thus, it suffices to ensure that the
 380 following two inequalities hold for $X \sim \mathcal{N}(0, 1)$:

$$\sup_{r \in \mathcal{P}(m), r \neq 0} \frac{r^2(y)}{\mathbf{E}[r^2(X)]} \leq \frac{1}{\alpha_y} \quad \text{and} \quad \sup_{q \in \mathcal{P}(m-1), q \neq 0} \frac{(B^2 - y^2)q^2(y)}{\mathbf{E}[(B^2 - X^2)q^2(X)]} \leq \frac{1}{\alpha_y}, \quad (1)$$

381 where we use Lemma 3.3 to show that $\mathbf{E}[(B^2 - X^2)q^2(X)] > 0$ for all non-zero polynomials
 382 $q \in \mathcal{P}(m-1)$. The first expression can be bounded using Lemma 3.4 when $m \leq C_1/\sqrt{\alpha}$.
 383 We now focus on the second expression. By Lemma 3.3, $\mathbf{E}_{X \sim \mathcal{N}(0,1)}[(B^2 - X^2)q^2(X)] \geq$
 384 $0.5 \mathbf{E}_{X \sim \mathcal{N}(0,1)}[B^2 q^2(X)]$. Therefore, we have that

$$\begin{aligned} \sup_{q \in \mathcal{P}(m-1), q \neq 0} \frac{(B^2 - y^2)q^2(y)}{\mathbf{E}_{X \sim \mathcal{N}(0,1)}[(B^2 - X^2)q^2(X)]} &\leq \sup_{q \in \mathcal{P}(m-1), q \neq 0} \frac{B^2 q^2(y)}{\mathbf{E}_{X \sim \mathcal{N}(0,1)}[(B^2 - X^2)q^2(X)]} \\ &\leq \sup_{q \in \mathcal{P}(m-1), q \neq 0} \frac{B^2 q^2(y)}{\mathbf{E}_{X \sim \mathcal{N}(0,1)}[0.5 B^2 q^2(X)]} = 2 \sup_{q \in \mathcal{P}(m-1), q \neq 0} \frac{q^2(y)}{\mathbf{E}_{X \sim \mathcal{N}(0,1)}[q^2(X)]}, \end{aligned}$$

385 where the first inequality uses that the denominator is positive and $y^2 q^2(y) \geq 0$ and the second
 386 inequality uses that $\mathbf{E}_{X \sim \mathcal{N}(0,1)}[(B^2 - X^2)q^2(X)] \geq 0.5 \mathbf{E}_{X \sim \mathcal{N}(0,1)}[B^2 q^2(X)]$. The expression
 387 above is of the same form as the first expression in Equation (1), and thus is also bounded above by
 388 $1/\alpha_y$ when $m \leq C_1/\sqrt{\alpha}$ using Lemma 3.4. This completes the proof of Theorem 2.5. \square

389 **Proof sketch of Lemma 3.3:** The proof of Lemma 3.3 is a relatively straightforward application
 390 of Hölder's inequality and the Gaussian Hypercontractivity Theorem (stated below). For $p \in (0, \infty)$,
 391 we define the L^p -norm of a random variable X to be $\|X\|_{L^p} := (\mathbf{E}[|X|^p])^{1/p}$.

392 **Fact 3.5** (Gaussian Hypercontractivity [Bog98, Nel73]). *Let $X \sim \mathcal{N}(0, 1)$. If $p \in \mathcal{P}(d)$ and $t \geq 2$,*
 393 *then $\|p(X)\|_{L^t} \leq (t-1)^{d/2} \|p(X)\|_{L^2}$.*

394 **Proof sketch of Lemma 3.4:** The proof is based on Hermite Analysis (see Appendix A for more
 395 details). The *normalized probabilist's Hermite polynomials*, $\{h_i, i \in [m]\}$ form a basis of $\mathcal{P}(m)$ and
 396 satisfy the property $\mathbf{E}_{X \sim \mathcal{N}(0,1)}[h_i(X)h_j(X)] = \mathbf{1}(i=j)$. Since r is a polynomial of degree at most
 397 m , we can represent $r(x) = \sum_{i=1}^m a_i h_i(x)$ for some $a_i \in \mathbb{R}$. Using orthonormality of h_i under the
 398 Gaussian measure, we get that $\mathbf{E}_{X \sim \mathcal{N}(0,1)}[r^2(X)] = \sum_{i=1}^m a_i^2$. By a standard optimization argument,
 399 we get that the supremum of $r^2(y)/\mathbf{E}[r^2(X)]$ is exactly $\sum_{i=1}^m h_i^2(y)$. It remains to show that for
 400 every $y \in \mathbb{R}$, $\sum_{i=1}^m \alpha_y h_i^2(y) \leq 1/2$. As $m \leq C/\sqrt{\alpha}$ for a small enough constant C , it suffices to
 401 show that for every $i \in [m]$, $\alpha_y h_i^2(y) \leq O(\sqrt{\alpha})$. As $\alpha_y := \sqrt{\alpha} \exp(-y^2(1-\alpha)/2)$, the following
 402 tail bound on the Hermite polynomials can be used:

403 **Lemma 3.6** ([Kra04]). *Let h_k be defined as above. Then $\max_{x \in \mathbb{R}} h_k^2(x) e^{-x^2/2} = O(k^{-1/6})$.*

404 We break our analysis in two cases:

405 **Case 1:** $|y| \leq 1/\sqrt{\alpha}$. Since $\alpha^2 y \leq 1$, Lemma C.3 implies that for every $|y| \leq 1/\sqrt{\alpha}$, $\alpha_y h_i^2(y) =$
 406 $\sqrt{\alpha} \exp(1) h_i^2(y) \exp(-y^2/2) = O(\sqrt{\alpha})$.

407 **Case 2:** $|y| > 1/\sqrt{\alpha}$. In this case, we use rather crude bounds. A direct calculation shows that
 408 $|h_i(x)| \leq i^i (1 + |x|)^i$. Since $\alpha \in (0, 1/2)$, we get that $\alpha_y h_i^2(y) \leq \sqrt{\alpha} \exp(-y^2/4 + i \log(2i|y|))$.
 409 It remains to show that $\exp(-y^2/4 + i \log(2i|y|)) = O(1)$ under given conditions on i and y . We
 410 have that $\exp(-y^2/4 + i \log(2i|y|)) = O(1)$ whenever $|y| = \Omega(\sqrt{i \log i})$. Since $|y| \geq 1/\sqrt{\alpha}$, the
 411 former condition is satisfied whenever $i = O(1/\sqrt{\alpha})$. This completes the proof sketch. \square

412 **References**

- 413 [AAR99] G. E. Andrews, R. Askey, and R. Roy. *Special Functions*. Encyclopedia of Mathematics
414 and its Applications. Cambridge University Press, 1999.
- 415 [BBH⁺20] M. Brennan, G. Bresler, S. B. Hopkins, J. Li, and T. Schramm. Statistical query
416 algorithms and low-degree tests are almost equivalent. *arXiv preprint arXiv:2009.06107*,
417 2020. To appear in *34th Annual Conference on Learning Theory (COLT 2021)*.
- 418 [BBV08] M. F. Balcan, A. Blum, and S. Vempala. A discriminative framework for clustering via
419 similarity functions. In *Proceedings of the 40th Annual ACM Symposium on Theory of
420 Computing*, pages 671–680, 2008.
- 421 [BFJ⁺94] A. Blum, M. Furst, J. Jackson, M. Kearns, Y. Mansour, and S. Rudich. Weakly learning
422 DNF and characterizing statistical query learning using Fourier analysis. In *Proceedings
423 of the Twenty-Sixth Annual Symposium on Theory of Computing*, pages 253–262, 1994.
- 424 [BNJT10] M. Barreno, B. Nelson, A. D. Joseph, and J. D. Tygar. The security of machine learning.
425 *Machine Learning*, 81(2):121–148, 2010.
- 426 [BNL12] B. Biggio, B. Nelson, and P. Laskov. Poisoning attacks against support vector machines.
427 In *Proceedings of the 29th International Conference on Machine Learning, ICML 2012*,
428 2012.
- 429 [Bog98] V. Bogachev. *Gaussian measures*. Mathematical surveys and monographs, vol. 62, 1998.
- 430 [CFJ13] T. Cai, J. Fan, and T. Jiang. Distributions of angles in random packing on spheres.
431 *Journal of Machine Learning Research*, 14(1):1837–1864, 2013.
- 432 [CLS20] S. Chen, J. Li, and Z. Song. Learning mixtures of linear regressions in subexponential
433 time via fourier moments. In *Proceedings of the 52nd Annual ACM SIGACT Symposium
434 on Theory of Computing*, pages 587–600, 2020.
- 435 [CSV17] M. Charikar, J. Steinhardt, and G. Valiant. Learning from untrusted data. In *Proceedings
436 of STOC 2017*, pages 47–60, 2017.
- 437 [DeV89] R. D. DeVeaux. Mixtures of linear regressions. *Computational Statistics & Data
438 Analysis*, 8(3):227–245, November 1989.
- 439 [Die01] T. E. Dielman. *Applied Regression Analysis for Business and Economics*.
440 Duxbury/Thomson Learning Pacific Grove, CA, 2001.
- 441 [DK19] I. Diakonikolas and D. M. Kane. Recent advances in algorithmic high-dimensional
442 robust statistics. *CoRR*, abs/1911.05911, 2019.
- 443 [DK20] I. Diakonikolas and D. M. Kane. Small covers for near-zero sets of polynomials and
444 learning latent variable models. In *Proceedings of the 61st Annual IEEE Symposium on
445 Foundations of Computer Science (FOCS 2020)*, 2020.
- 446 [DKK⁺16] I. Diakonikolas, G. Kamath, D. M. Kane, J. Li, A. Moitra, and A. Stewart. Robust
447 estimators in high dimensions without the computational intractability. In *Proceedings
448 of FOCS’16*, pages 655–664, 2016.
- 449 [DKK⁺19] I. Diakonikolas, G. Kamath, D. Kane, J. Li, J. Steinhardt, and A. Stewart. Sever: A robust
450 meta-algorithm for stochastic optimization. In *Proceedings of the 36th International
451 Conference on Machine Learning, ICML 2019*, pages 1596–1606, 2019.
- 452 [DKS17] I. Diakonikolas, D. M. Kane, and A. Stewart. Statistical query lower bounds for robust
453 estimation of high-dimensional gaussians and gaussian mixtures. In *58th IEEE Annual
454 Symposium on Foundations of Computer Science, FOCS 2017*, pages 73–84, 2017. Full
455 version at <http://arxiv.org/abs/1611.03473>.
- 456 [DKS18] I. Diakonikolas, D. M. Kane, and A. Stewart. List-decodable robust mean estimation
457 and learning mixtures of spherical gaussians. In *Proceedings of the 50th Annual ACM
458 SIGACT Symposium on Theory of Computing, STOC 2018*, pages 1047–1060, 2018. Full
459 version available at <https://arxiv.org/abs/1711.07211>.

- 460 [DKS19] I. Diakonikolas, W. Kong, and A. Stewart. Efficient algorithms and lower bounds for
461 robust linear regression. In *Proceedings of the Thirtieth Annual ACM-SIAM Symposium*
462 *on Discrete Algorithms, SODA 2019*, pages 2745–2754, 2019.
- 463 [EY07] A. Eremenko and P. Yuditskii. Uniform approximation of $\operatorname{sgn} x$ by polynomials and
464 entire functions. *Journal d’Analyse Mathématique*, 101(1):313–324, 2007.
- 465 [Fel16] V. Feldman. Statistical query learning. In *Encyclopedia of Algorithms*, pages 2090–2095.
466 Springer New York, 2016.
- 467 [Fel17] V. Feldman. A general characterization of the statistical query complexity. In *Proceedings*
468 *of the 30th Conference on Learning Theory, COLT 2017*, pages 785–830. PMLR, 2017.
- 469 [FGR⁺17] V. Feldman, E. Grigorescu, L. Reyzin, S. Vempala, and Y. Xiao. Statistical algorithms
470 and a lower bound for detecting planted cliques. *J. ACM*, 64(2):8:1–8:37, 2017.
- 471 [FGV17] V. Feldman, C. Guzman, and S. S. Vempala. Statistical query algorithms for mean vector
472 estimation and stochastic convex optimization. In *Proceedings of the Twenty-Eighth*
473 *Annual ACM-SIAM Symposium on Discrete Algorithms, SODA 2017*, pages 1265–1277.
474 SIAM, 2017.
- 475 [Gan02] M. I. Ganzburg. Limit theorems for polynomial approximation with hermite and freud
476 weights. *Approximation Theory X: Abstract and Classical Analysis (CK Chui, et al,*
477 *eds.)*, pages 211–221, 2002.
- 478 [GR08] M. I. Ganzburg and J. Rognes. *Limit theorems of polynomial approximation with*
479 *exponential weights*. American Mathematical Soc., 2008.
- 480 [HKP⁺17] S. B. Hopkins, P. K. Kothari, A. Potechin, P. Raghavendra, T. Schramm, and D. Steurer.
481 The power of sum-of-squares for detecting hidden structures. In *58th IEEE Annual*
482 *Symposium on Foundations of Computer Science, FOCS 2017*, pages 720–731. IEEE
483 Computer Society, 2017.
- 484 [Hop18] S. B. Hopkins. *Statistical inference and the sum of squares method*. PhD thesis, Cornell
485 University, 2018.
- 486 [HR09] P.J. Huber and E. M. Ronchetti. *Robust statistics*. Wiley New York, 2009.
- 487 [HRRS86] F. R. Hampel, E. M. Ronchetti, P. J. Rousseeuw, and W. A. Stahel. *Robust statistics. The*
488 *approach based on influence functions*. Wiley New York, 1986.
- 489 [HS17] S. B. Hopkins and D. Steurer. Efficient bayesian estimation from few samples: Commu-
490 nity detection and related problems. In *58th IEEE Annual Symposium on Foundations of*
491 *Computer Science, FOCS 2017*, pages 379–390. IEEE Computer Society, 2017.
- 492 [JJ94] M. I. Jordan and R. A. Jacobs. Hierarchical mixtures of experts and the EM algorithm.
493 *Neural Computation*, 6(2):181–214, 1994.
- 494 [KC20] J. Kwon and C. Caramanis. EM converges for a mixture of many linear regressions.
495 In *International Conference on Artificial Intelligence and Statistics*, pages 1727–1736.
496 PMLR, 2020.
- 497 [Kea98] M. J. Kearns. Efficient noise-tolerant learning from statistical queries. *Journal of the*
498 *ACM*, 45(6):983–1006, 1998.
- 499 [KKK19] S. Karmalkar, A. R. Klivans, and P. Kothari. List-decodable linear regression. In
500 *Advances in Neural Information Processing Systems 32: Annual Conference on Neural*
501 *Information Processing Systems 2019, NeurIPS 2019*, pages 7423–7432, 2019.
- 502 [Kra04] I. Krasikov. New bounds on the Hermite polynomials. *arXiv preprint math/0401310*,
503 2004.
- 504 [KS53] S. Karlin and L. S. Shapley. *Geometry of moment spaces*, volume 12. American
505 Mathematical Soc., 1953.

- 506 [LAT⁺08] J.Z. Li, D.M. Absher, H. Tang, A.M. Southwick, A.M. Casto, S. Ramachandran, H.M.
507 Cann, G.S. Barsh, M. Feldman, L.L. Cavalli-Sforza, and R.M. Myers. Worldwide human
508 relationships inferred from genome-wide patterns of variation. *Science*, 319:1100–1104,
509 2008.
- 510 [LL18] Y. Li and Y. Liang. Learning mixtures of linear regressions with nearly optimal com-
511 plexity. In *Conference On Learning Theory, COLT 2018*, pages 1125–1144. PMLR,
512 2018.
- 513 [LRV16] K. A. Lai, A. B. Rao, and S. Vempala. Agnostic estimation of mean and covariance. In
514 *Proceedings of FOCS’16*, 2016.
- 515 [McD09] J. H. McDonald. *Handbook of Biological Statistics, volume 2*. Sparky House Publishing,
516 Baltimore, MD, 2009.
- 517 [MV18] M. Meister and G. Valiant. A data prism: Semi-verified learning in the small-alpha
518 regime. In *Conference On Learning Theory, COLT 2018*, volume 75 of *Proceedings of*
519 *Machine Learning Research*, pages 1530–1546. PMLR, 2018.
- 520 [Nel73] E. Nelson. The free markoff field. *Journal of Functional Analysis*, 12(2):211–227, 1973.
- 521 [O’D14] R. O’Donnell. *Analysis of Boolean Functions*. Cambridge University Press, 2014.
- 522 [PLJD10] P. Paschou, J. Lewis, A. Javed, and P. Drineas. Ancestry informative markers for fine-
523 scale individual assignment to worldwide populations. *Journal of Medical Genetics*,
524 47:835–847, 2010.
- 525 [RL87] P. J. Rousseeuw and A. M. Leroy. *Robust Regression and Outlier Detection*. John Wiley
526 & Sons, Inc., New York, NY, USA, 1987.
- 527 [RPW⁺02] N. Rosenberg, J. Pritchard, J. Weber, H. Cann, K. Kidd, L.A. Zhivotovsky, and M.W.
528 Feldman. Genetic structure of human populations. *Science*, 298:2381–2385, 2002.
- 529 [RY20] P. Raghavendra and M. Yau. List decodable learning via sum of squares. In *Proceedings*
530 *of the 2020 ACM-SIAM Symposium on Discrete Algorithms, SODA 2020*, pages 161–180.
531 SIAM, 2020.
- 532 [SKL17] J. Steinhardt, P. W. Koh, and P. S. Liang. Certified defenses for data poisoning attacks.
533 In *Advances in Neural Information Processing Systems 30*, pages 3520–3532, 2017.
- 534 [SVC16] J. Steinhardt, G. Valiant, and M. Charikar. Avoiding imposters and delinquents: Adver-
535 sarial crowdsourcing and peer prediction. In *NIPS*, pages 4439–4447, 2016.
- 536 [Sze89] G. Szegő. *Orthogonal Polynomials*, volume XXIII of *American Mathematical Society*
537 *Colloquium Publications*. A.M.S, Providence, 1989.
- 538 [Val84] L. Valiant. A theory of the learnable. *Communications of the ACM*, 27(11):1134–1142,
539 1984.
- 540 [ZJD16] K. Zhong, P. Jain, and I. S. Dhillon. Mixed linear regression with multiple components.
541 In *Advances in Neural Information Processing Systems 29: Annual Conference on*
542 *Neural Information Processing Systems 2016*, pages 2190–2198, 2016.

543 Checklist

- 544 1. For all authors...
- 545 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s
546 contributions and scope? [Yes]
- 547 (b) Did you describe the limitations of your work? [Yes]
- 548 (c) Did you discuss any potential negative societal impacts of your work? [N/A]
- 549 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
550 them? [Yes]

- 551 2. If you are including theoretical results...
- 552 (a) Did you state the full set of assumptions of all theoretical results? [Yes]
- 553 (b) Did you include complete proofs of all theoretical results? [Yes]
- 554 3. If you ran experiments...
- 555 (a) Did you include the code, data, and instructions needed to reproduce the main experi-
- 556 mental results (either in the supplemental material or as a URL)? [N/A]
- 557 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
- 558 were chosen)? [N/A]
- 559 (c) Did you report error bars (e.g., with respect to the random seed after running experi-
- 560 ments multiple times)? [N/A]
- 561 (d) Did you include the total amount of compute and the type of resources used (e.g., type
- 562 of GPUs, internal cluster, or cloud provider)? [N/A]
- 563 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 564 (a) If your work uses existing assets, did you cite the creators? [N/A]
- 565 (b) Did you mention the license of the assets? [N/A]
- 566 (c) Did you include any new assets either in the supplemental material or as a URL? [N/A]
- 567
- 568 (d) Did you discuss whether and how consent was obtained from people whose data you're
- 569 using/curating? [N/A]
- 570 (e) Did you discuss whether the data you are using/curating contains personally identifiable
- 571 information or offensive content? [N/A]
- 572 5. If you used crowdsourcing or conducted research with human subjects...
- 573 (a) Did you include the full text of instructions given to participants and screenshots, if
- 574 applicable? [N/A]
- 575 (b) Did you describe any potential participant risks, with links to Institutional Review
- 576 Board (IRB) approvals, if applicable? [N/A]
- 577 (c) Did you include the estimated hourly wage paid to participants and the total amount
- 578 spent on participant compensation? [N/A]