L1 Regression with Lewis Weights Subsampling

Aditya Parulekar adityaup@cs.utexas.edu UT Austin Advait Parulekar advaitp@utexas.edu UT Austin Eric Price ecprice@cs.utexas.edu UT Austin

May 21, 2021

Abstract

We consider the problem of finding an approximate solution to ℓ_1 regression while only observing a small number of labels. Given an $n \times d$ unlabeled data matrix X, we must choose a small set of $m \ll n$ rows to observe the labels of, then output an estimate $\widehat{\beta}$ whose error on the original problem is within a $1 + \varepsilon$ factor of optimal. We show that sampling from X according to its Lewis weights and outputting the empirical minimizer succeeds with probability $1 - \delta$ for $m > O(\frac{1}{\varepsilon^2}d\log\frac{d}{\varepsilon\delta})$. This is analogous to the performance of sampling according to leverage scores for ℓ_2 regression, but with exponentially better dependence on δ . We also give a corresponding lower bound of $\Omega(\frac{d}{\varepsilon^2} + (d + \frac{1}{\varepsilon^2})\log\frac{1}{\delta})$.

1 Introduction

The standard linear regression problem is, given a data matrix $X \in \mathbb{R}^{n \times d}$ and corresponding values $y \in \mathbb{R}^n$, to find a vector $\beta \in \mathbb{R}^d$ minimizing $\|X\beta - y\|_p$. Least squares regression (p = 2) is the most common, but least absolute deviation regression (p = 1) is sometimes preferred for its robustness to outliers and heavy-tailed noise. In this paper we focus on ℓ_1 regression:

$$\beta^* = \underset{\beta \in \mathbb{R}^d}{\arg \min} \|X\beta - y\|_1 \tag{1}$$

But what happens if the unlabeled data X is cheap but the labels y are expensive? Can we choose a small subset of indices, only observe the corresponding labels, and still recover a good estimate $\hat{\beta}$ of the true solution? We would like an algorithm that works with probability $1 - \delta$ for any input (X, y); this necessitates that our choice of indices be randomized, so the adversary cannot concentrate the noise on them. Formally we define the problem as follows:

Problem 1 (Active L1 regression). There is a known matrix $X \in \mathbb{R}^{n \times d}$ and a fixed unknown vector y. A learner interacts with the instance by querying rows indexed $\{i_k\}_{k \in [m]}$ adaptively, and is shown labels $\{y_{i_k}\}_{k \in [m]}$ corresponding to the rows queried. The learner must return $\widehat{\beta}$ such that with probability $1 - \delta$ over the learner's randomness,

$$||X\widehat{\beta} - y||_1 \le (1 + \varepsilon) \min_{\beta} ||X\beta - y||_1.$$
 (2)

Some rows of X may be more important than others. For example, if one row is orthogonal to all the others, we need to query it to have any knowledge of the corresponding y; but if many rows are in the same direction it should suffice to label a few of them to predict the rest.

A natural approach to this problem is to attach some notion of "importance" p_1, \ldots, p_n to each row of X, then sample rows proportional to p_i . We can represent this as a "sampling-and-reweighting" sketch $S \in \mathbb{R}^{m \times n}$, where each row is $\frac{1}{p_i}e_i$ with probability proportional to p_i . This reweighting is such that $\mathbb{E}_S[\|Sv\|_1] \propto \|v\|_1$ for any vector v. By querying m rows we can observe Sy, and so can output the empirical risk minimizer (ERM)

$$\widehat{\beta} := \arg\min \|SX\beta - Sy\|_1. \tag{3}$$

For fixed β , $\mathbb{E}_S ||SX\beta - Sy||_1 \propto ||X\beta - y||_1$. The hope is that, if the p_i are chosen carefully, the ERM $\widehat{\beta}$ will satisfy (2) with relatively few samples. Our main result is that this is true if the p_i are drawn according to the ℓ_1 Lewis weights:

Theorem 1.1 (Informal). Problem 1 can be solved with $m = O(\frac{1}{\varepsilon^2}d\log\frac{d}{\varepsilon\delta})$ queries. For constant $\delta = \Theta(1)$, $m = O(\frac{1}{\varepsilon^2}d\log d)$ suffices.

Note that, while the model allows for adaptive queries, this algorithm is nonadaptive.

We next show that our sample complexity is near-optimal by demonstrating the following lower bound on the number of queries needed by any algorithm to obtain an accurate estimate.

Theorem 1.2 (Informal). Any algorithm satisfying Problem 1 must query $\Omega(d \log \frac{1}{\delta} + \frac{d}{\varepsilon^2} + \frac{1}{\varepsilon^2} \log \frac{1}{\delta})$ rows on some instances (X, y).

For small δ , the upper bound is the product of $d, \frac{1}{\varepsilon^2}$, and $\log(1/\delta)$ while the lower bound is the product of each pair.

1.1 Related Work

If all the labels are known: LAD regression cannot be solved in closed form. It can be written as a linear program, but this is relatively slow to solve. One approach to speeding up LAD regression is "sketch-and-solve," which replaces (1) with (3), which has fewer constraints and so can be solved faster. The key idea here is to acquire regression guarantees by ensuring that S is a subspace embedding for the column space of $[X \ y]$.

For a survey on techniques to do this, we direct the reader to [Woo14], [Mah11], [Cla05]. In [Woo14], the emphasis is on *oblivious* sketches – distributions which do not require knowledge of $[X \ y]$. On the other hand, [Mah11], [Cla05] discuss sketches that depend on $[X \ y]$. Most relevant to us [DLS18], which shows that sampling-and-reweighting matrices S using Lewis weights of $[X \ y]$ suffice; we give a simple proof of this in Remark 2.2. The problem is that figuring out which labels are important involves looking at all the labels.

Active ℓ_2 regression: Here we return to our setting, where only a subset of the labels is available to us. A number of works have studied this problem, including [DMM06, DW17, DM21]. The ℓ_2 version of the problem was solved optimally in [CP19], where an algorithm was given using $O(\frac{d}{\varepsilon})$ queries to find $\hat{\beta}$ satisfying $\mathbb{E}\left[\|X\hat{\beta}-y\|_2^2\right] \leq (1+\varepsilon)\|X\beta^*-y\|_2^2$. Independent, identical sampling using leverage scores achieves the same guarantee using $O(d\log d + \frac{d}{\varepsilon})$ queries. Note that these results for ℓ_2 ERM only work in expectation, while our results hold with high probability. One can get high probability bounds in the ℓ_2 setting by taking the median of $O(\log 1/\delta)$ repetitions, but the ERM itself does not succeed with high probability.

Subspace embedding for ℓ_1 norms: Subspace embeddings for the ℓ_1 norm have been studied in a long line of work including [Tal90], [Tal95], [LT89], [DDH⁺09], and [CP15], the most recent of which describes an iterative algorithm to approximate *Lewis weights*, which are the analogue of leverage scores for importance sampling preserving ℓ_1 norms. The [CP15] result shows that, for the same $m = O(\frac{1}{\varepsilon^2}d\log\frac{d}{\varepsilon\delta})$ sample complexity as given in Theorem 3.1, a sampler sketch S based on the Lewis weights of X will have $||SX\beta||_1 \approx_{\varepsilon} ||X\beta||_1$ for all $\beta \in \mathbb{R}^d$.

Our approach. At a very high level the goal of this paper is to replace the ℓ_2 leverage score analysis of the [CP19] active regression paper with the ℓ_1 Lewis weight analysis in the [CP15] subspace embedding paper. However, the differences between ℓ_1 and ℓ_2 are significant enough that very little of the [CP19] proof approach remains.

Per [CP15], the Lewis weight sampling-and-embedding matrix S preserves $||X\beta||_1$ for all β . The problem is that it doesn't preserve $||X\beta - y||_1$: if y has outliers, we have no idea where they are to sample them. In the ℓ_2 setting, this difficulty is addressed using the closed-form solution $\beta^* = X^{\dagger}y$. Then if S is a subspace embedding it will preserve $||X\beta - X\beta^*||$, so it suffices to bound the expectation of $||S(X\beta^* - y)||_2^2$. In the ℓ_1 setting, not only is β^* not expressible in closed form, but there can be many equally valid minimizers β^* that are far from each other. In Appendix A we show how this approach extends to the ℓ_1 setting to give a simple proof of Theorem 1.1 for a constant factor approximation (i.e., $\varepsilon = O(1)$); but the existence of multiple β^* makes $\varepsilon < 1$ seem unobtainable by this approach.

Instead, we massage the [CP19] subspace embedding proof into the appropriate form, as we discuss in Section 3. While S doesn't preserve the total error $||X\beta - y||_1$, it does preserve relative error $||X\beta - y||_1 - ||X\beta^* - y||_1$; the effect of outliers is canceled out, so that this concentrates similarly well to $||X\beta - X\beta^*||_1$. This approach would not work for ℓ_2 : the effect of outliers does not entirely cancel out there, since the square loss has unbounded influence.

Concurrent work: A very similar set of results appears concurrently and independently in [CD21]. Their main result is identical to ours, with a similar proof. They also extend the result to $1 , but with a significantly weaker <math>m = \widetilde{O}(d^2/\varepsilon^2)$ bound. They do not have the $\Omega(d \log \frac{1}{\delta})$ lower bound.

2 Preliminaries: Subspace Embeddings and Importance Sampling

A key idea used in our analysis is that of a ℓ_1 subspace embedding, which is a linear sketch of a matrix that preserves ℓ_1 norms within the column space of a matrix:

Definition 2.1 (Subspace Embeddings). A subspace embedding for the column space of the matrix $X \in \mathbb{R}^{n \times d}$ is a matrix S such that for all $\beta \in \mathbb{R}^d$,

$$||SX\beta|| = (1 \pm \varepsilon)||X\beta||$$

Remark 2.2. Consider the simpler setting in which we had access to all of y, but we still want to subsample rows to improve computational complexity. We can view the regression loss $||X\beta - y||_1$ as the ℓ_1 norm of the point $[X \ y] \begin{bmatrix} \beta \\ -1 \end{bmatrix}$ in the column space of $[X \ y]$. Indeed, suppose $\beta^* = \arg\min ||X\beta - y||_1$ as before and let $\widehat{\beta} = \arg\min ||SX\beta - Sy||_1$. Then, $\widehat{\beta}$ solves problem 1 because, for $\varepsilon < \frac{1}{3}$,

$$\|X\widehat{\beta} - y\|_1 \le \frac{1}{1-\varepsilon} \|SX\widehat{\beta} - Sy\|_1 \le \frac{1}{1-\varepsilon} \|SX\beta^* - Sy\|_1 \le \frac{1+\varepsilon}{1-\varepsilon} \|X\beta^* - y\|_1 \le (1+4\varepsilon) \|X\beta^* - y\|_1.$$

One way to construct a subspace embedding is by sampling rows and rescaling them appropriately:

Definition 2.3 (Sampling and Reweighting with $\{p_i\}_{i=1}^n$). For any sequence $\{p_i\}_{i=1}^n$, let $N = \sum_i p_i$. Then, the sampling-and-reweighting distribution $S(\{p_i\}_{i=1}^n)$ over the set of matrices $S \in \mathbb{R}^{N \times n}$ is such that each row of S is independently the ith standard basis vector with probability $\frac{p_i}{N}$, scaled by $\frac{1}{p_i}$. For any $k \in [N]$, let i_k denote the index such that $S_{k,i_k} = \frac{1}{p_{i_k}}$.

When working in ℓ_2 , there is a natural choice for re-weighting: the leverage scores of the rows [Woo14].

Definition 2.4 (Leverage Scores). The leverage score of the ith row of a matrix X, $l_i(X)$ is defined as $x_i^{\top}(X^{\top}X)^{-1}x_i$.

For ℓ_1 subspace embeddings, the analogous weights are the ℓ_1 Lewis weights, defined implicitly as the unique weights $\{w_i(X)\}_{i=1}^n$ that satisfy $w_i(X) = l_i(WX)$ where W is a diagonal matrix with ith diagonal entry $\frac{1}{\sqrt{w_i(X)}}$. We will drop the explicit dependence on X whenever it is clear from context.

Definition 2.5 (Lewis Weights). The ℓ_1 Lewis weights of a matrix X are the unique weights $\{w_i\}_{i=1}^n$ that satisfy $w_i^2 = x_i^\top (\sum_{j=1}^n \frac{1}{w_j} x_j x_j^\top)^{-1} x_i$ for all i.

Lewis weights are defined in general for general ℓ_p norms, but we will only need the ℓ_1 Lewis weights. For basic properties of Lewis weights, we direct the reader to [CP15]. Using these definitions, we now state the main consequence of using Lewis weights. This result comes from a line of work on embeddings from subspaces of $L_1[0,1]$ to ℓ_1^m such as [Tal90], but is reproduced here similar to how it is presented in [CP15].

Theorem 2.6 ([CP15] Theorem 2.3). Sampling at least $O(\frac{d \log d}{\varepsilon^2})$ rows according to the ℓ_1 Lewis weights $\{w_i\}_{i=1}^n$ of a matrix $X \in \mathbb{R}^{n \times d}$ results in a subspace embedding for X with at least some constant probability. If at least $O(\frac{d \log \frac{d}{\varepsilon \delta}}{\varepsilon^2})$ rows are sampled, then we have a subspace embedding with probability at least $1 - \delta$.

2.1 Properties of Lewis Weights

We will need some properties of Lewis weights, particularly of how they change when the matrix X is modified.

Lemma 2.7 ([CP15] Lemma 5.5). The ℓ_1 Lewis weights of a matrix do not increase when rows are added.

Lemma 2.8. Let $X \in \mathbb{R}^{n \times d}$, and let $X' \in \mathbb{R}^{kn \times d}$ be X stacked on itself k times, with each row scaled down by k. Then, each of the Lewis weights is reduced by a factor of k.

3 Proof Overview

Theorem 3.1. Let $X \in \mathbb{R}^{n \times d}$ have ℓ_1 Lewis weights $\{w_i\}_{i \in [n]}$, and let $0 < \varepsilon, \delta < 1$. Then, for any N that is at least $O\left(\frac{d}{\varepsilon^2}\log\frac{d}{\varepsilon\delta}\right)$, there is a sampling-and-reweighting distribution $\mathcal{S}(\{p_i\}_{i=1}^n)$ satisfying $\sum_i p_i = N$ such that for all y, if $S \sim \mathcal{S}(\{p_i\}_{i=1}^n)$ and $\widehat{\beta} = \arg\min \|SX\beta - Sy\|_1$, we have

$$||X\widehat{\beta} - y||_1 \le (1 + \varepsilon) \min_{\beta} ||X\beta - y||_1$$

with probability $1 - \delta$. If $\delta = O(1)$ is some constant, then N at least $O\left(\frac{1}{\varepsilon^2}d\log d\right)$ rows suffice.

Regression guarantees from column-space embeddings. As noted in Remark 2.2, it would suffice to show that $||SX\beta - Sy||_1 \approx ||X\beta - y||_1$ for all β . The problem is that this is impossible without knowing y: if one random entry of y is very large, we would need to sample it to estimate $||X\beta - y||_1$ accurately. However, we don't actually need to estimate $||X\beta - y||_1$; we just need to be able to distinguish values of β for which $||X\beta - y||_1$ is far from $||X\beta^* - y||_1$ from values for which it is close. That is, it would suffice to accurately

estimate
$$||X\hat{\beta} - y||_1 - ||X\beta^* - y||_1$$
 with $||SX\hat{\beta} - Sy||_1 - ||SX\beta^* - Sy||_1$ (4)

for every possible β . In the above example where y has a single large outlier coordinate, sampling this coordinate or not will dramatically affect both terms, but will not affect the difference very much. As such, our key lemma, Lemma 4.1, states that ℓ_1 Lewis weight sampling achieves (4) with high probability. In particular, using at least $m \geq O(\frac{d}{\varepsilon^2} \log \frac{d}{\varepsilon \delta})$ rows we have

$$(\|SX\beta^* - Sy\|_1 - \|SX\beta - Sy\|_1) - (\|X\beta^* - y\|_1 - \|X\beta - y\|_1) < \varepsilon \|X(\beta^* - \beta)\|_1$$
 (5)

for all β with probability at least $1-\delta$. We do this by adapting the argument of [CP15] which shows that S is a column-space embedding with high probability. We have summarized this argument below.

Column-space embedding using Lewis weights ([CP15]). An important result in [CP15], which directly implies the high probability subspace embedding, and which will be useful to us later is the following moment bound on deviations of $||SX\beta||_1$.

Lemma 3.2 ([CP15] Lemma 7.4). If N is at least $O\left(\frac{d}{\varepsilon^2}\log\frac{d}{\varepsilon\delta}\right)$, and $S \in \mathbb{R}^{N \times n}$ is drawn from the sampling-and-reweighting distribution $S(\{p_i\}_{i=1}^N)$ with $\sum_i p_i = N$ and $\{p_i\}_{i=1}^n$ proportional to Lewis weights $\{w_i\}_{i=1}^n$, then

$$\mathbb{E}_{S} \left[\left(\max_{\|X\beta\|_{1}=1} |\|SX\beta\|_{1} - \|X\beta\|_{1} | \right)^{l} \right] \leq \varepsilon^{l} \delta$$

The proof follows from this chain of inequalities:

$$\mathbb{E}\left[\left(\max_{\|X\beta\|_{1}=1}\|SX\beta\|_{1} - \|X\beta\|_{1}\right)^{l}\right] \stackrel{(A)}{\leq} 2^{l} \mathbb{E}\left[\left(\max_{\|X\beta\|_{1}=1}\left|\sum_{k}\sigma_{k}\frac{|x_{i_{k}}^{T}\beta|}{p_{i_{k}}}\right|\right)^{l}\right] \stackrel{(B)}{\leq} 2^{l} \mathbb{E}\left[\left(\max_{\|X\beta\|_{1}=1}\sum_{k}\sigma_{k}\frac{x_{i_{k}}^{T}\beta}{p_{i_{k}}}\right)^{l}\right] \stackrel{(C)}{\leq} \varepsilon^{l}\delta$$

where the σ_k are independent Rademacher variables, which are ± 1 with probability 1/2 each, and p_{i_k} is proportional to the ℓ_1 Lewis weight of row i_k . (A) follows by symmetrizing the objective $F := \max_{\|X\beta\|_1 = 1} \|SX\beta\|_1 - \|X\beta\|_1$. (B) follows from a contraction lemma. (C) is shown by constructing a related matrix with bounded Lewis weights and applying Lemma 4.5 from [Tal90] reproduced below.

Lemma 3.3. There exists constant C such that for any $X \in \mathbb{R}^{n \times d}$ with all ℓ_1 Lewis weights less than $C \frac{\varepsilon^2}{\log(\frac{n}{\lambda})}$ and $l = \log(2n/\delta)$, we have

$$\mathbb{E}_{\sigma} \left[\left(\max_{\|X\beta\|_1 = 1} \left| \sum_{i=1}^n \sigma_i x_i^{\top} \beta \right| \right)^l \right] \le \frac{\varepsilon^l \delta}{2}$$
 (6)

Regression Guarantees using Lewis weight sampling. In this work, we show the following chain of inequalities.

$$\mathbb{E}_{S} \left[\left(\max_{\|X\beta^{*} - X\beta\| = 1} \left| (\|SX\beta^{*} - Sy\|_{1} - \|SX\beta - Sy\|_{1}) - (\|X\beta^{*} - y\|_{1} - \|X\beta - y\|_{1}) \right| \right)^{l} \right] \\
\stackrel{(A)}{\leq} 2^{l} \mathbb{E}_{S,\sigma} \left[\left(\max_{\|X\beta^{*} - X\beta\| = 1} \left| \sum_{k} \sigma_{k} \left(\frac{|x_{i_{k}}^{\top} \beta^{*} - y_{i_{k}}|}{p_{i_{k}}} - \frac{|x_{i_{k}}^{\top} \beta - y_{i_{k}}|}{p_{i_{k}}} \right) \right| \right)^{l} \right] \\
\stackrel{(B)}{\leq} 2^{2l+1} \mathbb{E}_{S,\sigma} \left[\left(\max_{\|X(\beta^{*} - \beta)\|_{1} = 1} \left| \sum_{k} \sigma_{i_{k}} \frac{x_{i_{k}}^{\top}}{p_{i_{k}}} (\beta^{*} - \beta) \right| \right)^{l} \right]$$

$$\stackrel{(C)}{\leq} \varepsilon^{l} \delta$$

$$(7)$$

Here, for (A), we symmetrize the left hand side of (5) in Lemma 4.2. For (B), we apply a different contraction lemma, Lemma 4.3, that allows us to remove y from our expression, and then end up with the same moment bound for (C). Step (C) is essentially an application of Lemma 4.5 to SX, however, because we cannot immediately bound the Lewis weights of SX to confirm the constraints of the Lemma, we instead construct another matrix X'' which does not significantly alter the right hand side of inequality (7) while having bounded Lewis weights. This is done in Lemmas 4.6 and 4.7.

3.1 Lower Bounds

We will show that any algorithm must see $\Omega(d\log\frac{1}{\delta} + \frac{1}{\varepsilon^2}\log\frac{1}{\delta} + \frac{d}{\varepsilon^2})$ labels to return $\widehat{\beta}$ satisfying $\|X\widehat{\beta} - y\|_1 \le (1+\varepsilon)\|X\beta^* - y\|_1$ with probability greater than $1-\delta$.

For the lower bound proof it is convenient to consider a distributional version of the problem:

Problem 2 (Distributional active L1 regression). There is an unknown joint distribution P over a finite set $\mathcal{X} \times \mathcal{Y} \subset \mathbb{R}^d \times \mathbb{R}$, with $|\mathcal{Y}| = 2$. The learner is allowed to adaptively observe N i.i.d. samples from $P(\cdot|X=x)$ for the learner's choice of N values $x \in \mathcal{X}$. The learner must return $\widehat{\beta}$ satisfying

$$\mathbb{E}_{(X,Y)\sim P}\left[|X^{\top}\widehat{\beta} - Y|\right] \le (1+\varepsilon)\inf_{\beta} \mathbb{E}_{(X,Y)\sim P}\left[|X^{\top}\beta - Y|\right]. \tag{8}$$

with probability at least $1 - \delta$.

We begin with a lemma that shows that solving the original, Problem 1, for some n polynomial in the parameters d, ε, δ is harder than solving the distributional version, Problem 2.

Lemma 3.4. A randomized algorithm that solves Problem 1 for $n = \frac{2}{\varepsilon^2} \left(\log \frac{2}{\delta} + d \log \frac{3d}{\varepsilon} \right)$ with accuracy ε and failure probability δ can be used to solve any instance of Problem 2, where \mathcal{X}, \mathcal{Y} , in the unit ℓ_{∞} ball, with accuracy 6ε and failure probability 2δ , for small ε .

We then prove lower bounds on the accuracy for any algorithm on Problem 2.

In all our lower bounds, x is a uniform e_i , and $y \in \{0,1\}$. For $\Omega(\frac{d}{\varepsilon^2})$, we set $P(y|x=e_i)$ to $\frac{1}{2} \pm \varepsilon$ uniformly at random independently for each i; getting an ε -approximate solution requires getting most of the biases correct, which requires $\frac{1}{\varepsilon^2}$ samples from most of the coordinates e_i . The $\Omega(\frac{1}{\varepsilon^2}\log\frac{1}{\delta})$ instance sets $P(y|x=e_i)$ to $\frac{1}{2} \pm \varepsilon$ with the same bias for each i; solving this is essentially distinguishing a ε biased coin from a $-\varepsilon$ -biased coin. Finally, for $\Omega(d\log\frac{1}{\delta})$ we set $P(y|x=e_i)=0$ except for a random hidden i^* with $P(y \mid x=e_{i^*})=\frac{3}{4}$. Solving this instance requires finding i^* , but there's a δ chance the first $d\log\frac{1}{\delta}$ queries are all zero.

Theorem 3.5. For any $d \geq 2$, $\epsilon < \frac{1}{10}$, $\delta < \frac{1}{4}$, there exist sets $\mathcal{X} \in \mathbb{R}^d$, $\mathcal{Y} \in \mathbb{R}$ of inputs and labels, and a distribution P on $\mathcal{X} \times \mathcal{Y}$ such that any algorithm which solves Problem 2, with $\varepsilon = 1$, requires at least $m = \Omega(\frac{d}{\epsilon^2} + \frac{1}{\epsilon^2}\log\frac{1}{\delta} + d\log\frac{1}{\delta})$ samples.

4 Proof of Theorem 3.1

Lemma 4.1. Let $X \in \mathbb{R}^{n \times d}$ have ℓ_1 Lewis weights $\{w_i\}_{i \in [n]}$. Then, for any N that is at least $O\left(\frac{d}{\varepsilon^2}\log\frac{d}{\varepsilon\delta}\right)$, there is a sampling-and-reweighting distribution $\mathcal{S}(\{p_i\}_{i=1}^n)$ satisfying $\sum_i p_i = N$ such that for all y, if $S \sim \mathcal{S}(\{p_i\}_{i=1}^n)$ and $\beta^* = \arg\min \|X\beta - y\|_1$, we have for all β

$$(\|SX\beta^* - Sy\|_1 - \|SX\beta - Sy\|_1) - (\|X\beta^* - y\|_1 - \|X\beta - y\|_1) \le \varepsilon \cdot \|X\beta^* - X\beta\|_1$$
 (9)

with probability at least $1 - \delta$. Further, for constant δ , $m = O(d \log d/\varepsilon^2)$ rows suffice.

This lemma is proved for constant and high probability bounds in Section 4.1. Given this, we can prove the main theorem.

Proof of Theorem 3.1. Applying Lemma 4.1 to $\widehat{\beta} := \arg \min \|SX\beta - Sy\|_1$, we get

$$\left(\|SX\beta^* - Sy\|_1 - \|SX\widehat{\beta} - Sy\|_1\right) \le \left(\|X\beta^* - y\|_1 - \|X\widehat{\beta} - y\|_1\right) + \varepsilon \cdot \|X\beta^* - X\widehat{\beta}\|_1$$

Since $\widehat{\beta}$ is the minimizer of $||SX\beta - Sy||_1$, the left side is non-negative. So,

$$||X\widehat{\beta} - y||_1 \le ||X\beta^* - y||_1 + \varepsilon \cdot ||X\beta^* - X\widehat{\beta}||_1$$

$$\le ||X\beta^* - y||_1 + \varepsilon \cdot (||X\beta^* - y||_1 + ||X\widehat{\beta} - y||_1)$$

Rearranging, and assuming $\varepsilon < 1/2$,

$$||X\widehat{\beta} - y||_1 \le \frac{1+\varepsilon}{1-\varepsilon} ||X\beta^* - y||_1$$

$$\le (1+4\varepsilon)||X\beta^* - y||_1$$

Using $\varepsilon' = \varepsilon/4$ proves the theorem.

4.1 Proof of Lemma 4.1

This argument is similar to that in Appendix B of [CP15]. In order to prove Lemma 4.1, by Markov's inequality, it is sufficient to show that for some l,

$$M := \mathbb{E}_{S} \left[\left(\max_{\|X\beta^{*} - X\beta\| = 1} \left| (\|SX\beta^{*} - Sy\|_{1} - \|SX\beta - Sy\|_{1}) - (\|X\beta^{*} - y\|_{1} - \|X\beta - y\|_{1}) \right| \right] \leq \varepsilon^{l} \delta^{l} \delta^{l} \right]$$

To show this, we will symmetrize, then use a contraction lemma to cancel the y terms. Then, with all the terms being within the column space of SX, we use the fact that S is a subspace embedding with high probability. We present two different bounds, one used for the constant probability and one for the high probability cases, but the following intermediate bound is the same for the two:

Lemma 4.2. Given a matrix $X \in \mathbb{R}^{n \times d}$, let $S(\{p_i\}_{i \in [n]})$ be any sampling-and-reweighting disribution, and let i_k be the row-indices chosen by this sampling matrix such that $S_{k,i_k} = \frac{1}{p_{i_k}}$. Let σ_k be independent Rademacher variables that are ± 1 each with probability 0.5. Then,

$$M \le 2^l \underset{S,\sigma}{\mathbb{E}} \left[\left(\max_{\|X\beta^* - X\beta\| = 1} \left| \sum_k \sigma_k \left(\frac{|x_{i_k}^\top \beta^* - y_{i_k}|}{p_{i_k}} - \frac{|x_{i_k}^\top \beta - y_{i_k}|}{p_{i_k}} \right) \right| \right)^l \right]$$
 (10)

This is essentially standard symmetrization; the proof is in Appendix B. To simplify the expression and eliminate the terms involving the labels, we then use a theorem from [LT89]:

Lemma 4.3 ([LT89] Theorem 5). Let $\Phi : \mathbb{R}^+ \to \mathbb{R}^+$ be convex and increasing, and let $\phi_k : \mathbb{R} \to \mathbb{R}$ be contractions such that $\phi_k(0) = 0$ for all k. Let \mathcal{F} be a class of functions on $\{1, 2, 3, \ldots, n\}$, and $\|g(f)\|_{\mathcal{F}} = \sup_{f \in \mathcal{F}} |g(f)|$. Then,

$$\mathbb{E}_{\sigma}\left[\Phi\left(\frac{1}{2}\left\|\sum_{k}\sigma_{k}\phi_{k}(f(k))\right\|_{\mathcal{F}}\right)\right] \leq \frac{3}{2}\mathbb{E}_{\sigma}\left[\Phi\left(\left\|\sum_{k}\sigma_{k}f(k)\right\|_{\mathcal{F}}\right)\right]$$

Lemma 4.4. For any $y \in \mathbb{R}^n$, we have

$$\mathbb{E}_{S,\sigma} \left[\left(\max_{\|X\beta^* - X\beta\| = 1} \left| \sum_{k} \sigma_k \left(\frac{|x_{i_k}^\top \beta^* - y_{i_k}|}{p_{i_k}} - \frac{|x_{i_k}^\top \beta - y_{i_k}|}{p_{i_k}} \right) \right| \right)^l \right] \\
\leq 2^{l+1} \mathbb{E}_{S,\sigma} \left[\left(\max_{\|X\beta^* - X\beta\|_1 = 1} \left| \sum_{k} \sigma_k \left(\frac{x_{i_k}^\top \beta^* - x_{i_k}^\top \beta}{p_{i_k}} \right) \right| \right)^l \right] \tag{11}$$

Proof. We take $\Phi(x) = x^l$, which is convex and increasing for l > 1, let \mathcal{F} be the set of functions f_{β} where $f_{\beta}(k) = \frac{x_{i_k}^{\top} \beta^* - x_{i_k}^{\top} \beta}{p_{i_k}}$ and β satisfies $||X\beta^* - X\beta||_1 = 1$, and let ϕ_k be defined as

$$\phi_k(z) = \frac{|x_{i_k}^{\top} \beta^* - y_{i_k}|}{p_{i_k}} - \frac{|x_{i_k}^{\top} \beta^* - z p_{i_k} - y_{i_k}|}{p_{i_k}}$$

This satisfies

$$\phi_k(f_{\beta}(k)) = \phi_k\left(\frac{x_{i_k}^{\top}\beta^* - x_{i_k}^{\top}\beta}{p_{i_k}}\right) = \frac{|x_{i_k}^{\top}\beta^* - y_{i_k}|}{p_{i_k}} - \frac{|x_{i_k}^{\top}\beta - y_{i_k}|}{p_{i_k}}.$$

This is a contraction, since

$$|\phi_k(z_1) - \phi_k(z_2)| = \left| \frac{|x_{i_k}^\top \beta^* - z_2 p_{i_k} - y_{i_k}|}{p_{i_k}} - \frac{|x_{i_k}^\top \beta^* - z_1 p_{i_k} - y_{i_k}|}{p_{i_k}} \right|$$

$$\leq \frac{|z_1 p_{i_k} - z_2 p_{i_k}|}{p_{i_k}} \leq |z_1 - z_2|$$

Applying Lemma 4.3 with these parameters, we have

$$\mathbb{E}\left[\left(\frac{1}{2}\max_{\|X\beta^* - X\beta\| = 1}\left|\sum_{k} \sigma_{k}\left(\frac{|x_{i_{k}}^{\top}\beta^* - y_{i_{k}}|}{p_{i_{k}}} - \frac{|x_{i_{k}}^{\top}\beta - truey_{i_{k}}|}{p_{i_{k}}}\right)\right|\right)^{l}\right] \\
\leq \frac{3}{2}\mathbb{E}\left[\left(\max_{\|X\beta^* - X\beta\|_{1} = 1}\left|\sum_{k} \sigma_{k}\left(\frac{x_{i_{k}}^{\top}\beta^* - x_{i_{k}}^{\top}\beta}{p_{i_{k}}}\right)\right|\right)^{l}\right]$$

After taking the expectation with respect to S and multiplying both sides by 2^l , this gives the statement of the lemma.

From here, we use two separate results to show the appropriate row counts for the constant and high probability cases. The constant probability case is left for Appendix C.

For high probability row-counts, we use a lemma from [CP15]:

Lemma 4.5 (8.2, 8.3, 8.4 in [CP15]). There exists constant C such that for any $X \in \mathbb{R}^{n \times d}$ with all ℓ_1 Lewis weights less than $C \frac{\varepsilon^2}{\log(\frac{n}{\delta})}$ and $l = \log(2n/\delta)$, then

$$\mathbb{E}_{\sigma} \left[\left(\max_{\|X\beta\|_1 = 1} \left| \sum_{i=1}^n \sigma_i x_i^{\top} \beta \right| \right)^l \right] \le \frac{\varepsilon^l \delta}{2}$$
 (12)

We want a similar statement, but for arbitrary matrices, with no bounds placed on the Lewis weights. To do this, we construct a new, related matrix using the following lemma, which is proved in Appendix B:

Lemma 4.6 (Similar to [CP15] Lemma B.1). Let X be any matrix, and let W be the matrix that has the Lewis weights of X in the diagonal entries. Let $N \ge \frac{d}{\varepsilon^2} \log \frac{d}{\varepsilon \delta}$. There exist constants C_1, C_2, C_3 such that we can construct a matrix X' such that

- X' has C_1dN rows,
- $X'^{\top}W'^{-1}X' \succeq X^{\top}W^{-1}X$, (where W' is the matrix that has the Lewis weights of X' in the diagonal entries),
- $||X'\beta||_1 \le C_2 ||X\beta||_1$ for all β ,
- the Lewis weights of X' are bounded by $\frac{C_3}{N}$.

Lemma 4.7. Consider $X \in \mathbb{R}^{n \times d}$ with ℓ_1 Lewis weights w_i . Let p_i be some set of sampling values such that $N = \sum_i p_i$ and, for some constants C, C_1, C_4 ,

$$p_i \ge \frac{\log\left(\frac{N + C_1 N d}{\delta}\right)}{C\varepsilon^2} w_i$$

Then, if $N \ge C_4 \frac{d}{\varepsilon^2} \log \frac{d}{\varepsilon \delta}$ and if $S \sim \mathcal{S}(\{p_i\}_{i \in [n]})$, then

$$\mathbb{E}_{S,\sigma} \left[\left(\max_{\|X\beta\|_1 = 1} \left| \sum_{k=1}^{N} \sigma_k \frac{x_{i_k}^{\top} \beta}{p_{i_k}} \right| \right)^l \right] \le \frac{\varepsilon^l \delta}{2}$$
(13)

Proof of Lemma 4.7. Ideally the Lewis weights of SX would be bounded by $C\frac{\varepsilon^2}{\log \frac{N}{\delta}}$ and we could directly apply Lemma 4.5 to SX to obtain a bound on the moment. However, we do not know this. Instead, we first construct X' using X as described in Lemma 4.6. We then construct a new matrix X'' by stacking X' on top of SX. Define W'' to be the diagonal matrix consisting of the ℓ_1 Lewis weights of X''. Define, for convenience, $R = N + C_1Nd$, which is the number of rows X'' has.

We can bound the term on the left side of (13) by the same term, summing over the rows of X'' instead. That is,

$$\mathbb{E}_{S,\sigma} \left[\left(\max_{\|X\beta\|=1} \left| \sum_{k=1}^{N} \sigma_k \frac{x_{i_k}^{\top} \beta}{p_{i_k}} \right| \right)^l \right] \leq \mathbb{E}_{S,\sigma} \left[\left(\max_{\|X\beta\|=1} \left| \sum_{i=1}^{R} \sigma_i x_i''^{\top} \beta \right| \right)^l \right]$$

Our goal is to apply Lemma 4.5 to the right side. To do this, we need to show the correct bound on its Lewis weights, and then have the term be a maximum over $||X''\beta||_1 = 1$, rather than $||X\beta||_1 = 1$.

Bounding the Lewis weights of X''. By Lemma 2.7, the ℓ_1 Lewis weights of a matrix do not increase when more rows are added. So, the rows in X'' that are from X' have Lewis weights that are bounded above by $C_3 \frac{\varepsilon^2}{\log(\frac{d}{\varepsilon^\delta})}$. Further,

$$X''^{\top}W''^{-1}X'' = \sum_{i=1}^{R} \frac{1}{w_i''} x_i''(x_i'')^{\top}$$

$$\succeq \sum_{i=1}^{R-N} \frac{1}{w_k''} x_k''(x_k'')^{\top} \qquad \text{since } \sum_{i=kC_1 d^2+1}^{N} \frac{1}{w_i''} x_i''(x_i'')^{\top} \succeq 0$$

$$= X'^{\top}W'^{-1}X' \succ X^{\top}W^{-1}X.$$

So, any row $y_i = x_i/p_i$ in X'' that is from SX satisfies

$$\begin{aligned} w_i''^2 &= y_i^\top (X''^\top W''^{-1} X'')^{-1} y_i \leq y_i^\top (X^\top W^{-1} X)^{-1} y_i \\ &= \frac{1}{p_i^2} x_i^\top (X^\top W^{-1} X)^{-1} x_i \\ &\leq \left(\frac{C\varepsilon^2}{\log\left(\frac{R}{\delta}\right)} \frac{1}{w_i}\right)^2 \cdot w_i^2 = \left(\frac{C\varepsilon^2}{\log\left(\frac{R}{\delta}\right)}\right)^2 \end{aligned}$$

which means that all of the Lewis weights of X'' are less than the larger of $C \frac{\varepsilon^2}{\log(\frac{R}{\delta})}$ and $C_3 \frac{\varepsilon^2}{\log(\frac{d}{\varepsilon\delta})}$. Now, for small enough ε, δ , $\log \frac{R}{\delta} \leq \frac{C}{C_3} \log \frac{d}{\varepsilon\delta}$, we have the Lewis weight upper bound for all rows of X'' is $C \frac{\varepsilon^2}{\log(\frac{R}{\delta})}$.

Renormalizing to maximize over $||X''\beta||_1 = 1$: If we define the following

$$F := \max_{\|X\beta\|_1 = 1} |\|SX\beta\|_1 - \|X\beta\|_1|$$

then,

$$||X''\beta||_1 = ||SX\beta||_1 + ||X'\beta||_1 \le (1 + C_2 + F)||X\beta||_1$$

So, we get

$$\left(\max_{\|X\beta\|=1} \left| \sum_{k=1}^{R} \sigma_k x_k''^{\top} \beta \right| \right)^{l} \le (1 + C_2 + F)^{l} \left(\max_{\|X''\beta\|=1} \left| \sum_{k=1}^{R} \sigma_k x_k''^{\top} \beta \right| \right)^{l} \\
\le 2^{l-1} ((1 + C_2)^{l} + F^{l}) \left(\max_{\|X''\beta\|=1} \left| \sum_{k=1}^{R} \sigma_k x_k''^{\top} \beta \right| \right)^{l}$$

Taking expectations of either side over just the Rademacher variables,

$$\mathbb{E}\left[\left(\max_{\|X\beta\|=1}\left|\sum_{k=1}^{R}\sigma_k x_k''^{\top}\beta\right|\right)^l\right] \leq 2^{l-1}((1+C_2)^l + F^l) \mathbb{E}\left[\left(\max_{\|X''\beta\|=1}\left|\sum_{k=1}^{R}\sigma_k x_k''^{\top}\beta\right|\right)^l\right]$$

Applying Lemma 4.5 to X'': Since X'' has R rows, and the correct Lewis weight bound, we can simply apply Lemma 4.5 to the right side above

$$\mathbb{E}_{\sigma} \left[\left(\max_{\|X\beta\|=1} \left| \sum_{k=1}^{R} \sigma_k x_k''^{\top} \beta \right| \right)^l \right] \le 2^{l-1} ((1+C_2)^l + F^l)) \frac{\varepsilon^l \delta}{2}$$

Now, by Lemma 3.2, we know that $\mathbb{E}_S[F^l] \leq \varepsilon^l \delta$. So, taking the expectation with respect to the sampling matrices of either side of the above, we get, for small enough ε, δ ,

$$\mathbb{E}_{S,\sigma} \left[\left(\max_{\|X\beta\|=1} \left| \sum_{k=1}^{kC_1 d^2 + N} \sigma_k x_k''^{\top} \beta \right| \right)^l \right] \le 2^{l-1} ((1 + C_2)^l + \varepsilon^l \delta) \frac{\varepsilon^l \delta}{2} \le 2^l (1 + C_2)^l \frac{\varepsilon^l \delta}{2}$$

So, solving the problem for $\varepsilon' = \frac{\varepsilon}{2+2C_2}$ gives the correct bound.

Finally, we can show Lemma 4.1

Proof of Lemma 4.1. Take $l = \log(2n/\delta)$, $N = 5\frac{(1+C_1)C_3}{C}\frac{d}{\varepsilon^2}\log\frac{d}{\varepsilon^\delta}$. Then, we apply Lemma 4.2, Lemma 4.4, and Lemma 4.7 to get

$$M < 2^{2l} \varepsilon^l \delta$$

which, solving the problem for $\varepsilon/4$, gives the correct bound. Then, applying Markov's inequality, we get that with probability δ ,

$$\max_{\|X\beta^* - X\beta\| = 1} \left| (\|SX\beta^* - Sy\|_1 - \|SX\beta - Sy\|_1) - (\|X\beta^* - y\|_1 - \|X\beta - y\|_1) \right| \le \varepsilon$$

Finally, scaling up appropriately gives, in generality,

$$|(\|SX\beta^* - Sy\|_1 - \|SX\beta - Sy\|_1) - (\|X\beta^* - y\|_1 - \|X\beta - y\|_1)| \le \varepsilon \|X\beta^* - X\beta\|_1$$

References

- [CD21] Xue Chen and Michał Dereziński. Query complexity of least absolute deviation regression via robust uniform convergence. arXiv preprint arXiv:2102.02322, 2021.
- [Cla05] Kenneth L. Clarkson. Subgradient and sampling algorithms for l1 regression. In *Proceedings of the Sixteenth Annual ACM-SIAM Symposium on Discrete Algorithms*, SODA '05, page 257–266, USA, 2005. Society for Industrial and Applied Mathematics.
- [CP15] Michael B. Cohen and Richard Peng. Lp row sampling by lewis weights. In *Proceedings* of the Forty-Seventh Annual ACM Symposium on Theory of Computing, STOC '15, page 183–192, New York, NY, USA, 2015. Association for Computing Machinery.
- [CP19] Xue Chen and Eric Price. Active regression via linear-sample sparsification. In Alina Beygelzimer and Daniel Hsu, editors, *Proceedings of the Thirty-Second Conference on Learning Theory*, volume 99 of *Proceedings of Machine Learning Research*, pages 663–695, Phoenix, USA, 25–28 Jun 2019. PMLR.

- [DDH+09] Anirban Dasgupta, Petros Drineas, Boulos Harb, Ravi Kumar, and Michael W Mahoney. Sampling algorithms and coresets for \ell_p regression. SIAM Journal on Computing, 38(5):2060–2078, 2009.
- [DLS18] David Durfee, Kevin A. Lai, and Saurabh Sawlani. ℓ₁ regression using lewis weights preconditioning and stochastic gradient descent. In Sébastien Bubeck, Vianney Perchet, and Philippe Rigollet, editors, *Proceedings of the 31st Conference On Learning Theory*, volume 75 of *Proceedings of Machine Learning Research*, pages 1626–1656. PMLR, 06–09 Jul 2018.
- [DM21] Michał Derezinski and Michael W Mahoney. Determinantal point processes in randomized numerical linear algebra. *Notices of the American Mathematical Society*, 68(1), 2021.
- [DMM06] Petros Drineas, Michael W Mahoney, and Shan Muthukrishnan. Sampling algorithms for l 2 regression and applications. In *Proceedings of the seventeenth annual ACM-SIAM symposium on Discrete algorithm*, pages 1127–1136, 2006.
- [DW17] Michał Dereziński and Manfred K Warmuth. Unbiased estimates for linear regression via volume sampling. arXiv preprint arXiv:1705.06908, 2017.
- [Gil52] E. N. Gilbert. A comparison of signalling alphabets. *The Bell System Technical Journal*, 31(3):504–522, 1952.
- [LS20] T. Lattimore and C. Szepesvari. Bandit Algorithms. Cambridge University Press, 2020.
- [LT89] M. Ledoux and M. Talagrand. Comparison theorems, random geometry and some limit theorems for empirical processes. *Ann. Probab.*, 17(2):596–631, 04 1989.
- [Mah11] Michael W. Mahoney. Randomized algorithms for matrices and data. Found. Trends Mach. Learn., 3(2):123–224, February 2011.
- [Tal90] Michel Talagrand. Embedding subspaces of l1 into ln 1. Proceedings of the American Mathematical Society, 108(2):363–369, 1990.
- [Tal95] Michel Talagrand. Embedding subspaces of lp in lpn. In J. Lindenstrauss and V. Milman, editors, *Geometric Aspects of Functional Analysis*, pages 311–326, Basel, 1995. Birkhäuser Basel.
- [Woo14] David P. Woodruff. Sketching as a tool for numerical linear algebra. Found. Trends Theor. Comput. Sci., 10(1–2):1–157, October 2014.

A Constant-factor approximation

If we just want a constant factor approximation, we can take S to be a constant probability ℓ_1 -subspace embedding, so that $||X\beta||_1 \leq 2||SX\beta||_1$ with probability at least 0.9. We have

$$||X\widehat{\beta} - y||_{1} \leq ||X\widehat{\beta} - X\beta^{*}||_{1} + ||X\beta^{*} - y||_{1}$$

$$\leq 2||SX\widehat{\beta} - SX\beta^{*}||_{1} + ||X\beta^{*} - y||_{1}$$

$$\leq 2(||SX\widehat{\beta} - Sy||_{1} + ||SX\beta^{*} - Sy||_{1}) + ||X\beta^{*} - y||_{1}$$

$$\leq 4(||SX\beta^{*} - Sy||_{1}) + ||X\beta^{*} - y||_{1}$$

where in the last inequality, we have used the fact that $\widehat{\beta}$ is the minimizer of $||SX\beta - Sy||_1$. Now, by Markov's inequality, with probability 0.9, $||SX\beta^* - Sy||_1 \le 10||X\beta^* - y||_1$. So, we have with probability 0.81,

$$||X\widehat{\beta} - y||_1 \le 41||X\beta^* - y||_1$$

Since we only used a constant-factor subspace embedding, the row count would be $O(d \log d)$.

B Proofs of Lemmas

Lemma 2.8. Let $X \in \mathbb{R}^{n \times d}$, and let $X' \in \mathbb{R}^{kn \times d}$ be X stacked on itself k times, with each row scaled down by k. Then, each of the Lewis weights is reduced by a factor of k.

Proof. Let $\{w_i\}_{i=1}^n$ be the Lewis weights of X, and let $\{w_i'\}_{i=1}^{kn}$ be the Lewis weights of X'. Let x_i be the ith row of X, and similarly let x_i' be the ith row of X'. Let the ordering of the rows be such that $x_{jn+i}' = \frac{1}{k}x_i$ for $0 \le j < k$. Let W be the diagonal matrix where $W_{ii} = w_i$. Since Lewis weights are defined circularly, we just need to check that the suggested weights work, and by uniqueness, they will be correct.

We know that $w_i^2 = x_i^{\top} (X^{\top} W^{-1} X)^{-1} x_i$. Therefore, if we take W' to be the diagonal matrix of size $kn \times kn$, and set the diagonal entries to be the Lewis weights of X divided by k, repeated k times, then we have

$$X'^{\top}W'^{-1}X' = \sum_{i=1}^{kn} \frac{1}{w_i'} x_i' x_i'^{\top} = \sum_{i=1}^{kn} \frac{k}{w_i} x_i' x_i'^{\top} = k \sum_{i=1}^{n} \frac{k}{w_i} \cdot \frac{1}{k^2} x_i x_i^{\top}$$

In the last expression above, we are only summing over the first set of rows in X', which are the scaled rows of X, and then multiplying by k since they are repeated k times. Now,

$$k \sum_{i=1}^{n} \frac{k}{w_i} \cdot \frac{1}{k^2} x_i x_i^{\top} = \sum_{i=1}^{n} \frac{1}{w_i} x_i x_i^{\top} = X^{\top} W^{-1} X$$

So, finally, for an arbitrary row x'_{jn+i} , which corresponds to row x_i in the original matrix, we get its Lewis weight:

$$w_{jn+i}^{\prime 2} = x_{jn+i}^{\prime \top} (X^{\prime \top} W^{\prime - 1} X^{\prime})^{-1} x_{jn+i}^{\prime} = \frac{1}{k^2} x_i^{\top} (X^{\top} W^{-1} X)^{-1} x_i = \frac{w_i^2}{k^2}$$

which proves that our suggested Lewis weights are consistent.

Lemma 4.2. Given a matrix $X \in \mathbb{R}^{n \times d}$, let $S(\{p_i\}_{i \in [n]})$ be any sampling-and-reweighting disribution, and let i_k be the row-indices chosen by this sampling matrix such that $S_{k,i_k} = \frac{1}{p_{i_k}}$. Let σ_k be independent Rademacher variables that are ± 1 each with probability 0.5. Then,

$$M \le 2^l \underset{S,\sigma}{\mathbb{E}} \left[\left(\max_{\|X\beta^* - X\beta\| = 1} \left| \sum_k \sigma_k \left(\frac{|x_{i_k}^\top \beta^* - y_{i_k}|}{p_{i_k}} - \frac{|x_{i_k}^\top \beta - y_{i_k}|}{p_{i_k}} \right) \right| \right)^l \right]$$

$$(10)$$

Proof. We proceed by symmetrization. Since the matrix S scales the rows by the probability they are picked with, the expectation of $||SM\beta||_1$ is just $||M\beta||_1$, for any matrix M and vector β . So, adding or subtracting the same term with a different sampling matrix S', $(||S'X\beta^* - S'y||_1 - ||S'X\beta - S'y||_1) - (||X\beta^* - y||_1 - ||X\beta - y||_1)$, is just adding a mean zero term, and since taking the lth power of a maximum is convex, this can only increase the expectation. That is,

$$\mathbb{E}_{S,S'} \left[\left(\max_{\|X\beta^* - X\beta\| = 1} | (\|SX\beta^* - Sy\|_1 - \|SX\beta - Sy\|_1) - (\|X\beta^* - y\|_1 - \|X\beta - y\|_1) | \right)^{l} \right] \\
\leq \mathbb{E}_{S,S'} \left[\left(\max_{\|X\beta^* - X\beta\| = 1} | ((\|SX\beta^* - Sy\|_1 - \|SX\beta - Sy\|_1) - (\|X\beta^* - y\|_1 - \|X\beta - y\|_1)) \\
- \left((\|S'X\beta^* - S'y\|_1 - \|S'X\beta - S'y\|_1) - (\|X\beta^* - y\|_1 - \|X\beta - y\|_1) | \right)^{l} \right] \\$$

So, we can bound M as

$$M \leq \mathbb{E}_{S,S'} \left[\left(\max_{\|X\beta^* - X\beta\| = 1} |\left(\|SX\beta^* - Sy\|_1 - \|SX\beta - Sy\|_1 \right) - \left(\|S'X\beta^* - S'y\|_1 - \|S'X\beta - S'y\|_1 \right) |\right)^{l} \right]$$

Let i_k be the indices chosen by S, and i'_k the indices chosen by S'. Rewriting this as a sum,

$$M \leq \underset{S,S'}{\mathbb{E}} \left[\left(\max_{\|X\beta^* - X\beta\| = 1} \left| \sum_{k} \left(\frac{|x_{i_k}^{\top} \beta^* - y_{i_k}|}{p_{i_k}} - \frac{|x_{i_k}^{\top} \beta - y_{i_k}|}{p_{i_k}} \right) - \sum_{k} \left(\frac{|x_{i_k}^{\top} \beta^* - y_{i_k}|}{p_{i_k'}} - \frac{|x_{i_k'}^{\top} \beta - y_{i_k'}|}{p_{i_k'}} \right) \right| \right)^{l} \right]$$

Now, since i_k and i'_k are independent and identically distributed, randomly swapping elements from either sum does not change the distribution. This amounts to adding a random sign σ_k to the terms, where $\sigma_k = \pm 1$ independently with probability 1/2. So,

$$\begin{split} M &\leq \underset{S,S',\sigma}{\mathbb{E}} \left[\left(\max_{\|X\beta^* - X\beta\| = 1} \left| \sum_{k} \sigma_k \left(\frac{|x_{i_k}^\top \beta^* - y_{i_k}|}{p_{i_k}} - \frac{|x_{i_k}^\top \beta - y_{i_k}|}{p_{i_k}} \right) - \right. \right. \\ &\left. \sum_{K} \sigma_k \left(\frac{|x_{i_k'}^\top \beta^* - y_{i_k'}|}{p_{i_k'}} - \frac{|x_{i_k'}^\top \beta - y_{i_k'}|}{p_{i_k'}} \right) \right| \right)^l \right] \\ &\leq \underset{S,S',\sigma}{\mathbb{E}} \left[\left(\max_{\|X\beta^* - X\beta\| = 1} \left| \sum_{k} \sigma_k \left(\frac{|x_{i_k}^\top \beta^* - y_{i_k}|}{p_{i_k}} - \frac{|x_{i_k}^\top \beta - y_{i_k}|}{p_{i_k}} \right) \right| + \right. \right. \\ &\left. \max_{\|X\beta^* - X\beta\| = 1} \left| \sum_{k} \sigma_k \left(\frac{|x_{i_k'}^\top \beta^* - y_{i_k'}|}{p_{i_k'}} - \frac{|x_{i_k'}^\top \beta - y_{i_k'}|}{p_{i_k'}} \right) \right| \right)^l \right] \\ &\leq 2^l \underset{S,\sigma}{\mathbb{E}} \left[\left(\max_{\|X\beta^* - X\beta\| = 1} \left| \sum_{k} \sigma_k \left(\frac{|x_{i_k}^\top \beta^* - y_{i_k}|}{p_{i_k}} - \frac{|x_{i_k}^\top \beta - y_{i_k}|}{p_{i_k}} \right) \right| \right)^l \right] \end{split}$$

Where the final inequality follows from $(a + b)^l \le 2^{l-1}(a^l + b^l)$. Putting these together,

$$M \le 2^{l} \underset{S,\sigma}{\mathbb{E}} \left[\left(\max_{\|X\beta^{*} - X\beta\| = 1} \left| \sum_{k} \sigma_{k} \left(\frac{|x_{i_{k}}^{\top}\beta^{*} - y_{i_{k}}|}{p_{i_{k}}} - \frac{|x_{i_{k}}^{\top}\beta - y_{i_{k}}|}{p_{i_{k}}} \right) \right| \right)^{l} \right]$$

$$(14)$$

Lemma 4.6 (Similar to [CP15] Lemma B.1). Let X be any matrix, and let W be the matrix that has the Lewis weights of X in the diagonal entries. Let $N \ge \frac{d}{\varepsilon^2} \log \frac{d}{\varepsilon \delta}$. There exist constants C_1, C_2, C_3 such that we can construct a matrix X' such that

- X' has C_1dN rows,
- $X'^{\top}W'^{-1}X' \succeq X^{\top}W^{-1}X$, (where W' is the matrix that has the Lewis weights of X' in the diagonal entries),
- $||X'\beta||_1 < C_2 ||X\beta||_1$ for all β ,
- the Lewis weights of X' are bounded by $\frac{C_3}{N}$.

Proof. Given matrix X, we can use Lemma B.1 from [CP15] to construct a new matrix X_1 that satisfies

- X_1 has C_1d^2 rows,
- $X_1^{\top}W_1^{-1}X_1 \succeq X^{\top}W^{-1}X$, (where W_1 is the matrix that has the Lewis weights of X_1 in the diagonal entries),
- $||X_1\beta||_1 \le C_2 ||X_1\beta||_1$ for all β ,
- the Lewis weights of X_1 are bounded by $\frac{C_3}{d}$.

So, we can take this matrix and stack it on itself $k = \frac{N}{d}$ times, while scaling each row down by the same k. This will be our matrix X'. X' will then have $k = C_1Nd$ rows, which satisfies the first

bullet. Also, by Lemma 2.8, this shrinks the Lewis weights by a factor of k, which changes the Lewis weight upper bound to

$$\frac{C_3}{kd} = \frac{C_3}{N}$$

which is what we need. Now, since we are repeating rows k times, but each row is scaled down by k, we have $||X_1\beta||_1 = ||X'\beta||_1$ for all β . Therefore, $||X'\beta||_1 \le C_2||X\beta||_1$ for all β . Finally, as in the proof of Lemma 2.8, we know that since we have duplicated the rows of X_1 k times but scaled them down by k, $X_1^\top W_1^{-1} X_1 = X'^\top W'^{-1} X'$, and so we are done.

C Proof for constant failure probability

For the constant probability row-count, we use a lemma from [LT89]:

Lemma C.1 ([LT89]). There exists a constant C such that for any matrix X with all Lewis weights less than $C \frac{\varepsilon^2}{\log d}$,

$$\mathbb{E} \left[\max_{\boldsymbol{\sigma}} \sum_{\boldsymbol{\parallel} \boldsymbol{X} \boldsymbol{\beta} \parallel_{1} = 1} \sum_{k} \sigma_{k} x_{i}^{\top} \boldsymbol{\beta} \right] \leq \varepsilon$$

In [LT89], this is proved with absolute values within the sum (that is, summing $\sigma_i|x_i^{\top}\beta|$). However, the first step of the proof removes these absolute values using a comparison lemma, bounding the term with absolute values by twice the term without absolute values.

Lemma C.2. For matrix X with ℓ_1 Lewis weights w_i , let p_i be some set of sampling values such that $\sum_i p_i = N$ and $p_i \ge \frac{\log d}{C\varepsilon^2} w_i$. If you sample $S \sim \mathcal{S}(\{p_i\}_{i \in [n]})$, then

$$\mathbb{E}_{S,\sigma} \left[\max_{\|X\beta\|_1 = 1} \left| \sum_{k} \sigma_k \frac{x_{i_k}^{\top} \beta}{p_{i_k}} \right| \right] \le \varepsilon \tag{15}$$

Proof. This proof is very similar to that of Lemma 4.7.

Construct X' using X as described in Lemma 4.6, with $N = \frac{C_3 \log d}{C}$. We then construct a new matrix X'' by stacking X' on top of SX. Define W'' to be the diagonal matrix consisting of the ℓ_1 Lewis weights of X''.

We can bound the term on the left side of (15) by the same term, summing over the rows of X'' instead. That is,

$$\mathbb{E}_{S,\sigma} \left[\max_{\|X\beta\|=1} \left| \sum_{k=1}^{N} \sigma_k \frac{x_{i_k}^{\top} \beta}{p_{i_k}} \right| \right] \leq \mathbb{E}_{S,\sigma} \left[\max_{\|X\beta\|=1} \left| \sum_{i=1}^{R} \sigma_i x_i''^{\top} \beta \right| \right]$$

Our goal is to apply Lemma C.1 to the right side. To do this, we need to show the correct bound on its Lewis weights, and then have the term be a maximum over $||X''\beta||_1 = 1$, rather than $||X\beta||_1 = 1$.

Bounding the Lewis weights of X''. By Lemma 2.7, the ℓ_1 Lewis weights of a matrix do not increase when more rows are added. So, the rows in X'' that are from X' have Lewis weights that

are bounded above by $\frac{C\varepsilon^2}{\log d}$. Further,

$$X''^{\top}W''^{-1}X'' = \sum_{i=1}^{R} \frac{1}{w_i''} x_i''(x_i'')^{\top}$$

$$\succeq \sum_{i=1}^{R-N} \frac{1}{w_k''} x_k''(x_k'')^{\top} \qquad \text{since } \sum_{i=kC_1 d^2+1}^{N} \frac{1}{w_i''} x_i''(x_i'')^{\top} \succeq 0$$

$$= X'^{\top}W'^{-1}X' \succ X^{\top}W^{-1}X.$$

So, any row $y_i = x_i/p_i$ in X'' that is from SX satisfies

$$\begin{split} w_i''^2 &= y_i^\top (X''^\top W''^{-1} X'')^{-1} y_i \leq y_i^\top (X^\top W^{-1} X)^{-1} y_i \\ &= \frac{1}{p_i^2} x_i^\top (X^\top W^{-1} X)^{-1} x_i \\ &\leq \left(\frac{C\varepsilon^2}{\log d} \frac{1}{w_i}\right)^2 \cdot w_i^2 = \left(\frac{C\varepsilon^2}{\log d}\right)^2 \end{split}$$

which means that all of the Lewis weights of X" are less than $\frac{C\varepsilon^2}{\log d}$.

Renormalizing to maximize over $||X''\beta||_1 = 1$: If we define the following

$$F \coloneqq \max_{\|X\beta\|_1 = 1} |\|SX\beta\|_1 - \|X\beta\|_1|$$

then,

$$||X''\beta||_1 = ||SX\beta||_1 + ||X'\beta||_1 \le (1 + C_2 + F)||X\beta||_1$$

So, we get

$$\max_{\|X\beta\|=1} \left| \sum_{k=1}^{R} \sigma_k x_k''^\top \beta \right| \le (1 + C_2 + F) \cdot \max_{\|X''\beta\|=1} \left| \sum_{k=1}^{R} \sigma_k x_k''^\top \beta \right|$$

Taking expectations of either side over just the Rademacher variables,

$$\mathbb{E}\left[\max_{\sigma}\left[\max_{\|X\beta\|=1}\left|\sum_{k=1}^{R}\sigma_{k}x_{k}^{"\top}\beta\right|\right] \leq (1+C_{2}+F)\mathbb{E}\left[\max_{\sigma}\left[\max_{\|X''\beta\|=1}\left|\sum_{k=1}^{R}\sigma_{k}x_{k}^{"\top}\beta\right|\right]\right]$$

Applying Lemma C.1 to X'': Since X'' has R rows, and the correct Lewis weight bound, we can simply apply Lemma C.1 to the right side above

$$\mathbb{E}_{\sigma} \left[\max_{\|X\beta\|=1} \left| \sum_{k=1}^{R} \sigma_k x_k''^{\top} \beta \right| \right] \le (1 + C_2 + F) \varepsilon$$

Now, by Lemma 3.2, we know that $\mathbb{E}_S[F] \leq \varepsilon$. So, taking the expectation with respect to the sampling matrices of either side of the above, we get, for small enough ε ,

$$\mathbb{E}_{S,\sigma} \left[\max_{\|X\beta\|=1} \left| \sum_{k=1}^{kC_1 d^2 + N} \sigma_k x_k''^{\top} \beta \right| \right] \le 2(1 + C_2)\varepsilon$$

So, solving the problem for $\varepsilon' = \frac{\varepsilon}{2+2C_2}$ gives the correct bound.

Therefore, we can similarly prove the constant-probability case for Lemma 4.1:

Proof of 4.1 for constant probability. We take $l=1,\ N=\frac{d}{\varepsilon^2}\log d$ and apply Lemma 4.2, Lemma 4.4, and Lemma C.2.

D Lower Bounds

We prove three main theorems that allow us to show Theorem 3.5: Theorems D.1, D.2, and D.4. To do this, we make several Claims, which are proved in section 7.1. Recall the reduction between the matrix problem and the distribution:

Lemma 3.4. A randomized algorithm that solves Problem 1 for $n = \frac{2}{\varepsilon^2} \left(\log \frac{2}{\delta} + d \log \frac{3d}{\varepsilon} \right)$ with accuracy ε and failure probability δ can be used to solve any instance of Problem 2, where \mathcal{X}, \mathcal{Y} , in the unit ℓ_{∞} ball, with accuracy 6ε and failure probability 2δ , for small ε .

Proof. Let $n = \frac{8}{\varepsilon^2} \left(\log \frac{2}{\delta} + d \log \frac{4d}{\varepsilon} \right)$. Construct an instance of Problem 1 in which the rows of feature matrix **X** and the corresponding label vector y are drawn i.i.d. from P. Let H be the unit ℓ_{∞} ball. We have the following:

Claim D.1. For all $\beta \in H$, with probability at least $1 - \delta$,

$$(1 - \varepsilon) \mathbb{E}_{(X,Y) \sim P} \left[|X^{\top} \beta - Y| \right] \le \frac{1}{n} ||\mathbf{X}\beta - y||_1 \le (1 + \varepsilon) \mathbb{E}_{(X,Y) \sim P} \left[|X^{\top} \beta - Y| \right]$$

Let β° denote the minimizer $\inf_{\beta} \mathbb{E}_{(X,Y)\sim P}\left[|X^{\top}\beta - Y|\right]$. Let β^{*} denote the minimizer of the matrix instance $\inf_{\beta} \|\mathbf{X}\beta - y\|_{1}$, and let $\widehat{\beta}$ denote the output of the algorithm on the instance generated. Then we have

$$(1 - \varepsilon) \mathbb{E}_{(X,Y) \sim P} \left[|X^{\top} \widehat{\beta} - Y| \right] \leq \frac{1}{n} \|\mathbf{X} \widehat{\beta} - y\|_{1}$$

$$\leq (1 + \varepsilon) \frac{1}{n} \|\mathbf{X} \beta^{*} - y\|_{1} \quad \text{with probability } 1 - \delta$$

$$\leq (1 + \varepsilon) \frac{1}{n} \|\mathbf{X} \beta^{\circ} - y\|_{1}$$

$$\leq (1 + \varepsilon)^{2} \mathbb{E}_{(X,Y) \sim P} \left[|X^{\top} \beta^{\circ} - Y| \right]$$

So with probability $1-2\delta$,

$$\mathbb{E}_{(X,Y)\sim P}\left[|X^{\top}\widehat{\beta} - Y|\right] \le (1 + 6\varepsilon) \,\mathbb{E}_{(X,Y)\sim P}\left[|X^{\top}\beta^{\circ} - Y|\right].$$

Theorem D.1. For any $d \geq 2$ and $\varepsilon < \frac{1}{10}$, there exist families $\mathcal{X} \in \mathbb{R}^d$, $\mathcal{Y} \in \mathbb{R}$ of inputs and labels respectively such that any algorithm which solves Problem 2 with $\delta < \frac{1}{4}$ requires at least $m = \frac{3d}{2000\varepsilon^2}$ samples.

We take \mathcal{X} to be the set of standard basis vectors, and the distribution over \mathcal{X} to be uniform. We will define a set \mathcal{B} as being a subset of the unit hypercube $\{-1,1\}^d$ such that every element is sufficiently far from every other.

Claim D.2. There is a set $\mathcal{B} \subset \mathcal{H}$ with $|\mathcal{B}| \geq 2^{0.2d}$ such that for any two $\beta_1, \beta_2 \in \mathcal{B}$, we have $|\beta_1 - \beta_2| > 0.2d$

Proof. Here we just need an error correcting code with constant rate and constant relative (Hamming) distance. The existence of such a code follows from the Gilbert-Varshamov bound [Gil52]. \Box

Fix some unknown β^* . We will have $Y = ZX^{\top}\beta^*$ where Z is an independent random variable with probability $\frac{1}{2} + \varepsilon$ of being 1, and $\frac{1}{2} - \varepsilon$ of being -1. This completes our description of P. We define $l(\beta)$ to be the ℓ_1 norm of the residuals for β , that is, $l(\beta) = \mathbb{E}_{(X,Y)\sim P}[|X^{\top}\beta - Y|]$. We have the following properties of $l(\beta)$.

Claim D.3. For D, \mathcal{B} as chosen above, $l(\beta^*) = 1 - 2\varepsilon$.

Claim D.4. For D, \mathcal{B} as chosen above, we have for all $\beta \in \mathcal{B}$, $l(\beta) - l(\beta^*) = \frac{2\varepsilon}{d} ||\beta - \beta^*||_1$.

Proof of Theorem D.1. Suppose some algorithm returns $\widehat{\beta}$ with $l(\widehat{\beta}) < (1 + \frac{\varepsilon}{5})l(\beta^*) \implies ||\beta^* - \widehat{\beta}||_1 < 0.1d$ with probability $\frac{3}{4}$. By Fano's inequality,

$$H(\beta^*|\widehat{\beta}) < H\left(\frac{1}{4}\right) + \frac{\log|\mathcal{B}| - 1}{4} < 0.05d,$$

and we have a lower bound on the mutual information between the output of our algorithm and the true parameter: $I(\hat{\beta}; \beta^*) = H(\beta^*) - H(\beta^*|\hat{\beta}) \ge 0.15d$. For an upper bound on the mutual information after seeing m samples, we use the data processing inequality.

$$I(\beta^*; \widehat{\beta}) \leq I(\beta^*; (Y_i)_{i \in [m]}) \leq \sum_{i=1}^m I(\beta^*; Y_i | (Y_j)_{j \in [i-1]})$$

$$= \sum_{i=1}^m H(Y_i | (Y_j)_{j \in [i-1]}) - H(Y_i | \beta^*, (Y_j)_{j \in [i-1]})$$

$$\leq \sum_{i=1}^m 1 - H(Y_i | \beta^*, I_i)$$

$$\leq 4\varepsilon^2 m$$

Here we have used that

$$H(Y_i|\beta^*, (Y_j)_{j \in [i-1]}) \ge H(Y_i|\beta^*, I_i, (Y_j)_{j \in [i-1]})$$

= $H(Y_i|\beta^*, I_i)$

and that the distribution of Y_i conditioned on β^* , I_i is just an independent Bernoulli with parameter $\frac{1}{2} + \varepsilon$ and so

$$\sum_{i=1}^{m} 1 - H(Y_i | \beta^*, I_i) \le \sum_{i=1}^{m} \left[1 + \left(\frac{1}{2} + \varepsilon \right) \log \left(\frac{1}{2} + \varepsilon \right) + \left(\frac{1}{2} - \varepsilon \right) \log \left(\frac{1}{2} - \varepsilon \right) \right] \le 4\varepsilon^2 m$$

So $0.15d \leq I(\beta^*; \widehat{\beta}) \leq 4\varepsilon^2 m$, and so we need $m \geq \frac{3d}{80\varepsilon^2}$. The result follows by replacing ε with 5ε .

We can use the same instance to give a high probability lower bound of $\Omega(\log \frac{1}{\delta}/\varepsilon^2)$.

Theorem D.2. For any d and $\varepsilon < \frac{1}{10}$, there exist sets $\mathcal{X} \in \mathbb{R}$, $\mathcal{Y} \in \mathbb{R}$ of inputs and labels respectively, and a distribution P on $\mathcal{X} \times \mathcal{Y}$ such that any algorithm which solves problem 2 requires at least $m = \frac{1}{4\varepsilon^2} \log \frac{1}{\delta}$ samples.

Proof. Consider two instances, denoted by subscripts (1) and (2) with $\beta_{(1)}^* = -\mathbb{1}_d$ and $\beta_{(2)}^* = \mathbb{1}_d$, where $\mathbb{1}_d \in \mathbb{R}^d$ is the all-ones vector. Denote by $P_{(i)}$ the distribution over \mathcal{X}, \mathcal{Y} for instance (i), and let $l_{\beta_{(i)}^*}(\beta) = \mathbb{E}_{(X,Y) \sim P_{(i)}}[|X^\top \beta - Y|]$ for $i \in \{1,2\}$.

Claim D.5. For any
$$\beta$$
, $\max\{\ell_{\beta_{(1)}^*}(\beta) - \ell_{\beta_{(1)}^*}(\beta_{(1)}^*), \ell_{\beta_{(2)}^*}(\beta) - \ell_{\beta_{(2)}^*}(\beta_{(2)}^*)\} > 2\varepsilon$

From this claim together with Claim D.3, we have for some $i \in \{1, 2\}$, $l_{\beta_{(i)}^*}(\beta) \ge (1+2\varepsilon)l_{\beta_{(i)}^*}(\beta_{(i)}^*)$, for all β .

Denote by $\widehat{\beta}$ the output of the algorithm. Denote by $\mathbb{P}_{(1)}$ the distribution over outputs by a algorithm interacting instance (1), and by $\mathbb{P}_{(2)}$ the distribution over outputs by a algorithm interacting instance (2). Denote by A the event that $\ell_{\beta_{(1)}^*}(\widehat{\beta}) - \ell_{\beta_{(1)}^*}(\beta_{(1)}^*) \geq 2\varepsilon$. Note that under A^c , we have $\ell_{\beta_{(2)}^*}(\widehat{\beta}) - \ell_{\beta_{(2)}^*}(\beta_{(2)}^*) \geq 2\varepsilon$. Because the algorithm fails with probability at most δ on any instance, we have $2\delta \geq \mathbb{P}_{(1)}(A) + \mathbb{P}_{(2)}(A^c)$. On the other hand, $\mathbb{P}_{(1)}(A) + \mathbb{P}_{(2)}(A^c) \geq e^{-D(\mathbb{P}_{(1)}||\mathbb{P}_{(2)})}$. We can bound the KL-divergence of the two distributions as an aggregate KL-divergence over the course of acquiring the samples.

Theorem D.3 (Lemma 15.1, [LS20]). If a learner interacts with two environments (1) and (2) through a policy $\pi(\cdot|I_1, Y_1, I_2, Y_2, \cdots, Y_{i-1})$ which dictates a distribution over actions I_i conditioned on the past $(I_1, Y_1, \cdots, Y_{i-1})$, and sees label Y_i distributed according to some label distribution $P_{(1),I_i}$ and $P_{(2),I_i}$, then the KL-divergence between the output of the learner on instance (1) and (2), $\mathbb{P}_{(1)}$ and $\mathbb{P}_{(2)}$ is given by

$$D(\mathbb{P}_{(1)}||\mathbb{P}_{(2)}) = \sum_{k=1}^{d} \mathbb{E}_{(1)} \left[\sum_{i=1}^{N} \mathbb{1}\{I_i = k\} \cdot D(P_{(1),I_i}||P_{(2),I_i}) \right]$$

Now, $P_{(1),k}$ is a Bernoulli with parameter $\frac{1}{2} + \varepsilon$, and $P_{(1),k}$ is a Bernoulli with parameter $\frac{1}{2} - \varepsilon$, so $D(P_{(1),k}||P_{(1),k}) \le 16\varepsilon^2$, and so we have

$$\begin{split} \sum_{k=1}^{d} \mathbb{E}_{(1)} \left[\sum_{i=1}^{N} \mathbb{1}\{I_{i} = k\} \cdot D(P_{(1),I_{i}} || P_{(2),I_{i}}) \right] &\leq \sum_{k=1}^{d} \mathbb{E}_{(1)} \left[\sum_{i=1}^{N} \mathbb{1}\{I_{i} = k\} \cdot 16\varepsilon^{2} \right] \\ &= 16\varepsilon^{2} \cdot \mathbb{E}_{(1)} \left[\sum_{k=1}^{d} \sum_{i=1}^{N} \mathbb{1}\{I_{i} = k\} \right] = 16\varepsilon^{2} m \end{split}$$

Putting this together, we have $\delta \geq e^{-16\varepsilon^2 m} \implies m \geq \frac{1}{16\varepsilon^2} \log \frac{1}{\delta}$, and the result follows by replacing ε with $\frac{1}{2}\varepsilon$.

Theorem D.4. For any $d \geq 2$, there exist sets $\mathcal{X} \in \mathbb{R}^d$, $\mathcal{Y} \in \mathbb{R}$ of inputs and labels, and a distribution P on $\mathcal{X} \times \mathcal{Y}$ such that any algorithm which solves Problem 2, with $\varepsilon = 1$, requires at least $m = \frac{d}{3} \log \frac{1}{8\delta}$ samples.

Proof. All logarithms are base 4. Consider instances in which $\mathcal{X} = \{e_1, e_2, \cdots, e_d\}$ where e_i denotes the *i*th standard basis vector and the distribution over \mathcal{X} is uniform. We take $Y = ZX^{\top}\beta^*$ for some β^* , where Z is an independent Bernoulli random variable which is 1 with probability $\frac{3}{4}$, and 0 otherwise. Consider d instances labelled with subscripts $(1), (2), \cdots, (d)$, one in which each of the d standard basis is β^* , that is, $\beta^*_{(i)} = e_i$. Denote by β_j the jth coordinate of β . For each instance, we have

Claim D.6. For all
$$i \in [d], \beta \in \mathbb{R}^d$$
, we have $\ell_{\beta_{(i)}^*}(\beta) \geq \frac{1}{4d}$ with equality when $\beta = \beta_{(i)}^*$

We would like our algorithm to return an estimate $\widehat{\beta}$ which satisfies $\ell_{\beta^*}(\widehat{\beta}) < \frac{1}{2d}$. We first note that any choice of β only succeeds to be this close to the optimal on a single instance.

Claim D.7. Any
$$\beta \in \mathbb{R}^d$$
 can only satisfy $\ell_{\beta_{(i)}^*}(\widehat{\beta}) < \frac{1}{2d}$ for one $i \in [d]$.

So, we may as well enforce that the algorithm return one of e_1, e_2, \dots, e_d , since any other output can be mapped to one of these to improve the performance of the algorithm.

We will allow our algorithm to sample $N = \frac{d}{3} \log \frac{1}{\delta}$ rows total. Let \mathcal{E} be the event that $Y_1, Y_2, \ldots Y_N$ are all zero. Given any algorithm \mathcal{A} , let $F_{\mathcal{A}}$ denote the set of rows it samples fewer than $\log \frac{1}{\delta}$ times with probability at least $\frac{1}{2}$, in event \mathcal{E} . Because the total number of rows sampled is $\frac{d}{3} \log \frac{1}{\delta}$, there must be at least $\frac{2d}{3}$ rows which are sampled fewer than $\frac{1}{2} \log \frac{1}{\delta}$ times in expectation.

By Markov's inequality, these rows are sampled fewer than $\log \frac{1}{\delta}$ times with probability at least $\frac{1}{2}$, and are thus all in $F_{\mathcal{A}}$. Let $B_{\mathcal{A}}$ denote the distribution over outputs $\widehat{\beta}$ of \mathcal{A} in event \mathcal{E} . Let $i_{\mathcal{A}} = \arg \min_{j \in F_{\mathcal{A}}} B_{\mathcal{A}}(j)$. Denote by $G_{\mathcal{A}}$ the event that row $i_{\mathcal{A}}$ is sampled fewer than $\log \frac{1}{\delta}$ times; by construction we have $\mathbb{P}(G_{\mathcal{A}}) > \frac{1}{2}$.

The subscripts are explicit because $F_{\mathcal{A}}, B_{\mathcal{A}}, i_{\mathcal{A}}, \mathbb{P}[G_{\mathcal{A}}]$ are properties of the algorithm and are independent of the instance with which it interacts. Consider the performance of this algorithm against the instance $\beta_{(i_{\mathcal{A}})}^*$.

Let $Y_{(i_{\mathcal{A}}),j,k}$ denote the label returned to the algorithm when it queries e_j for the kth time. Let $T_{(i_{\mathcal{A}})} = \min\{t | Y_{(i_{\mathcal{A}}),i_{\mathcal{A}},t} = 1\}$. Denote by $E_{(i_{\mathcal{A}})}$ the event that $T_{(i_{\mathcal{A}})} > \log \frac{1}{\delta}$. Because $T_{(i_{\mathcal{A}})}$ is a geometric random variable, we have $\mathbb{P}[E_{(i_{\mathcal{A}})}] > \delta$.

Now condition on the event $G_{\mathcal{A}} \cap E_{i_{\mathcal{A}}}$, which is an event with probability $\frac{1}{2}\delta$. Here our algorithm samples $i_{\mathcal{A}}$ fewer than $T_{i_{\mathcal{A}}}$ times, so it never sees a 1 and its output distribution is $B_{\mathcal{A}}$. It returns $i \in F_{\mathcal{A}} \setminus \{i_{\mathcal{A}}\}$ with probability at least $1 - B_{\mathcal{A}}(i_{\mathcal{A}}) \geq 1 - \frac{1}{|F_{\mathcal{A}}|} \geq 1 - \frac{3}{2d} \geq \frac{1}{4}$. In summary, even after $\frac{d}{3}\log\frac{1}{\delta}$ queries, no algorithm can return $\widehat{\beta}$ with $\|X\widehat{\beta} - y\| < (1+\varepsilon)\|X\beta^* - y\|$ with probability greater than $\frac{1}{8}\delta$. The result follows by replacing δ by 8δ .

Corollary D.5. Any algorithm that solves Problem 1 takes at least $\Omega(d \log \frac{1}{\delta} + \frac{d}{\varepsilon^2} + \frac{1}{\varepsilon^2} \log \frac{1}{\delta})$ samples for some $n = O(\frac{d \log \frac{d}{\delta}}{\varepsilon})$.

Proof. Each of the instances that demonstrate the lower bounds above, in Lemmas D.1, D.2, and D.4, take $|\mathcal{X}| = d$, the results follows from Lemma 3.4.

D.1 Proof of Claims D.1, D.3, D.4, D.6, and D.7

Claim D.1. For all $\beta \in H$, with probability at least $1 - \delta$,

$$(1 - \varepsilon) \mathbb{E}_{(X,Y) \sim P} \left[|X^{\top} \beta - Y| \right] \le \frac{1}{n} ||\mathbf{X}\beta - y||_1 \le (1 + \varepsilon) \mathbb{E}_{(X,Y) \sim P} \left[|X^{\top} \beta - Y| \right]$$

Proof of Claim D.1. By assumption, we know that $X^{\top}\beta, Y \in [-1, 1]$, so, $|X^{\top}\beta - Y| \in [0, 2]$. So, for fixed β , by Hoeffding's on the rows of $\mathbf{X}\beta - y$, we have that if $n \geq \frac{8}{\varepsilon^2} \log \frac{2}{\delta'}$, then with probability at least $1 - \delta'$,

$$\left(1 - \frac{\varepsilon}{2}\right) \mathbb{E}_{(X,Y) \sim P} \left[|X^{\top}\beta - Y| \right] \le \frac{1}{n} ||\mathbf{X}\beta - y||_{1} \le \left(1 + \frac{\varepsilon}{2}\right) \mathbb{E}_{(X,Y) \sim P} \left[|X^{\top}\beta - Y| \right] \tag{16}$$

Now, we construct a $\frac{\varepsilon}{2d}$ -covering S of the unit ℓ_{∞} ball H, with fewer than $\left(\frac{4d}{\varepsilon}\right)^d$ elements, so that for any β , there is some $\beta_c \in S$ such that $\|\beta - \beta_c\|_{\infty} \leq \frac{\varepsilon}{2d}$. To do this, simply take $S = \{\beta : \beta_i = k\frac{\varepsilon}{2d}, k \in \mathbb{Z} \cap [-2d/\varepsilon, 2d/\varepsilon]\}$.

Note that **X** has rows on the hypercube. So, if we denote $x_{i,j}$ to be the entry of **X** in the *i*th row and *j*th column, then $x_{i,j} \in \{-1,1\}$. Therefore, for any β ,

$$\|\mathbf{X}\beta\|_1 = \sum_{i=1}^n |x_i^\top \beta| \le \sum_{i=1}^n \sum_{j=1}^d |x_{i,j}\beta_j| \le \sum_{i=1}^n \sum_{j=1}^d |\beta_j| \le nd\|\beta\|_{\infty}$$

Therefore, we can apply Hoeffding's, as in (16), with $\delta' = \delta \left(\frac{\varepsilon}{4d}\right)^d$, and union bound over the set S, to get that for any $\beta \in S$, with probability at least $1 - \delta$, (16) holds.

Then, for any $\beta \in H$, by the covering property, we can find some $\beta_c \in S$ such that

$$\|\beta - \beta_c\|_{\infty} \le \frac{\varepsilon}{d} \implies \|\mathbf{X}\beta - \mathbf{X}\beta_c\|_1 \le n\varepsilon.$$
 (17)

We have

$$\|\mathbf{X}\beta_c - y\|_1 - \|\mathbf{X}\beta_c - \mathbf{X}\beta\|_1 \le \|\mathbf{X}\beta - y\|_1 \le \|\mathbf{X}\beta - \mathbf{X}\beta_c\|_1 + \|\mathbf{X}\beta_c - y\|_1$$

So, combining (16) and (17), and dividing by n, we finally have that if $n \geq \frac{8}{\varepsilon^2} \left(\log \frac{2}{\delta} + d \log \frac{4d}{\varepsilon} \right)$, then for all $\beta \in H$,

$$(1 - \varepsilon) \mathbb{E}_{(X,Y) \sim P} \left[|X^{\top} \beta - Y| \right] \le \frac{1}{n} ||\mathbf{X}\beta - y||_1 \le (1 + \varepsilon) \mathbb{E}_{(X,Y) \sim P} \left[|X^{\top} \beta - Y| \right]$$

Claim D.3. For D, \mathcal{B} as chosen above, $l(\beta^*) = 1 - 2\varepsilon$.

Proof of Claim D.3. The ℓ_1 error for the correct β is given by

$$\mathbb{E}_{(X,Y)\sim P} \left| X^{\top} \beta^* - Y \right|$$

$$= \mathbb{E}_X \left[E_{Y\sim P(\cdot|X)} \middle| |X^{\top} \beta^* - Y \middle| \right] \qquad \text{by independence}$$

$$= \mathbb{E}_X \left[\left(\frac{1}{2} + \varepsilon \right) \middle| X^{\top} \beta^* - X^{\top} \beta^* \middle| + \left(\frac{1}{2} - \varepsilon \right) \middle| X^{\top} \beta^* + X^{\top} \beta^* \middle| \right]$$

$$= \mathbb{E}_X \left[(1 - 2\varepsilon) \middle| X^{\top} \beta^* \middle| \right] \qquad \beta^* \in \mathcal{H}$$

$$= 1 - 2\varepsilon$$

Claim D.4. For D, \mathcal{B} as chosen above, we have for all $\beta \in \mathcal{B}$, $l(\beta) - l(\beta^*) = \frac{2\varepsilon}{d}||\beta - \beta^*||_1$.

Proof of Claim D.4.

$$\begin{split} &\mathbb{E}_{(X,Y)\sim P} \left| X^{\top}\beta - Y \right| \right| \\ &= \mathbb{E}_{X} \left[E_{Y\sim P(\cdot|X)} \middle| X^{\top}\beta - Y \right| \middle| \right] \\ &= \mathbb{E}_{X} \left[\left(\frac{1}{2} + \varepsilon \right) \middle| X^{\top}\beta - X^{\top}\beta^{*} \middle| + \left(\frac{1}{2} - \varepsilon \right) \middle| X^{\top}\beta + X^{\top}\beta^{*} \middle| \right] \\ &= (1 - 2\varepsilon) + 2\varepsilon \, \mathbb{E}_{X} [X^{\top}\beta - X^{\top}\beta^{*}] \\ &= (1 - 2\varepsilon) + 2\varepsilon \frac{1}{d} ||\beta - \beta^{*}||_{1} \end{split}$$

Claim D.5. For any β , $\max\{\ell_{\beta_{(1)}^*}(\beta) - \ell_{\beta_{(1)}^*}(\beta_{(1)}^*), \ell_{\beta_{(2)}^*}(\beta) - \ell_{\beta_{(2)}^*}(\beta_{(2)}^*)\} > 2\varepsilon$ Proof of Claim D.5.

$$l(\beta) + l(\beta) = 2 - 4\varepsilon + \frac{2\varepsilon}{d} \|\beta_{(1)}^* - \beta\|_1 + \frac{2\varepsilon}{d} \|\beta_{(2)}^* - \beta\|_1$$
$$\geq 2 - 4\varepsilon + \frac{2\varepsilon}{d} \|\beta_{(2)}^* - \beta_{(1)}^*\|_1$$
$$= 2$$

$$\implies \max\{\ell_{\beta_{(1)}^*}(\beta) - \ell_{\beta_{(1)}^*}(\beta_{(1)}^*), \ell_{\beta_{(2)}^*}(\beta) - \ell_{\beta_{(2)}^*}(\beta_{(2)}^*)\} > 2\varepsilon, \forall \beta \in \mathbb{R}^d$$

Claim D.6. For all $i \in [d], \beta \in \mathbb{R}^d$, we have $\ell_{\beta_{(i)}^*}(\beta) \geq \frac{1}{4d}$ with equality when $\beta = \beta_{(i)}^*$ Proof of Claim D.6.

$$\ell_{\beta_{(i)}^*}(\beta) = \frac{1}{d} \sum_{j \neq i} |\beta_j| + \frac{\frac{1}{2} + \varepsilon}{d} |1 - \beta_i| + \frac{\frac{1}{2} - \varepsilon}{d} |\beta_i|$$
$$\geq \frac{\frac{1}{2} - \varepsilon}{d} (|\beta_i| + |1 - \beta_i|) + \frac{2\varepsilon}{d} |1 - \beta_i| \geq \frac{\frac{1}{2} - \varepsilon}{d}$$

Claim D.7. Any $\beta \in \mathbb{R}^d$ can only satisfy $\ell_{\beta_{(i)}^*}(\widehat{\beta}) < \frac{1}{2d}$ for one $i \in [d]$.

Proof of Claim D.7. Indeed, suppose β was such that $\ell_{\beta_{(I)}^*}(\beta), \ell_{\beta_{(J)}^*}(\beta) < \frac{1}{2d}$. Then we must have

$$\frac{1}{2d} \ge \ell_{\beta_{(I)}^*}(\beta)$$

$$= \frac{1}{d} \sum_{j \ne I} |\beta_j| + \frac{\frac{1}{2} - \varepsilon}{d} (|\beta_I| + |1 - \beta_i|) + \frac{2\varepsilon}{d} |1 - \beta_I|$$

$$\ge \frac{1}{d} \sum_{j \ne I} |\beta_j| + \frac{\frac{1}{2} - \varepsilon}{d} + \frac{2\varepsilon}{d} |1 - \beta_I|$$

$$\iff \varepsilon \ge \sum_{j \ne I} |\beta_j| + 2\varepsilon |1 - \beta_I|$$

$$\ge \sum_{j \ne I} |\beta_j| + 2\varepsilon - 2\varepsilon |\beta_I|$$

$$\iff 2|\beta_I| \ge ||\beta||_1 + 2\varepsilon$$

Similarly for J, so we would have $\|\beta\| \ge |\beta_I| + |\beta_J| \ge \|\beta\|_1 + 2\varepsilon$.