

# Conditional Linear Regression

Diego Calderon\*  
University of Arkansas  
dacalder@uark.edu

Brendan Juba† Sirui Li Zongyi Li  
Washington University in St. Louis  
{bjuba, sirui.li, zli}@wustl.edu

Lisa Ruan ‡  
M.I.T.  
llruan@mit.edu

## Abstract

Work in machine learning and statistics commonly focuses on building models that capture the vast majority of data, possibly ignoring a segment of the population as outliers. However, there does not often exist a good model on the whole dataset, so we seek to find a small subset where there exists a useful model. We are interested in finding a linear rule capable of achieving more accurate predictions for just a segment of the population. We give an efficient algorithm with theoretical analysis for the conditional linear regression task, which is the joint task of identifying a significant segment of the population, described by a  $k$ -DNF, along with its linear regression fit.

## 1 Introduction

Linear regression is the task of modeling the relationship between a result variable and some explanatory variables by a linear rule. Linear regression is a standard tool of statistical analysis, with widespread applications spanning essentially all of the sciences. While the standard linear regression task seeks to model the majority of the data, we consider problems where a regression fit could exist for some subset or portion of the data, that does not necessarily model the majority of the data. We will consider cases in which the subset with a linear model is described by some simple condition; in other words, we desire to perform linear regression on this conditional distribution. Note that neither the condition nor the model is known in advance. We are seeking an algorithm for this task that scales reasonably with the number of predictors used in the model, the dimension of the data, and so on.

To illustrate our problem, consider a set of patient data from a hospital that includes multiple continuous factors, such as rate of smoking, radiation exposure, physical activity. Assume we want to predict risk of developing lung cancer. There may be no linear rules that model, e.g., the risk of developing lung cancer for the majority of the data. However, there may be some linear model that fits a specific subset of the data well, such as adult city-dwellers. If such a model exists, we aim to find it together with the description of the corresponding subset. Our focus is on identifying the portions of the population for which such simply structured models succeed at making accurate predictions, even when such models do not exist for most of the population.

---

\*Part of this work was performed during an REU at Washington University in St. Louis, supported by WUSEF.

†Supported by an AFOSR Young Investigator Award and NSF award CCF-718380.

‡Part of this work was performed during an REU at Washington University in St. Louis, supported by the NSF Big Data Analytics REU Site, award IIS-1560191.

## 1.1 Our contributions

This problem was introduced by Juba (2017), who gave an algorithm for conditional linear regression under the  $\ell_\infty$ -loss where the predictor factors are sparse (i.e., its time and data requirements are exponential in the number of regression factors), and an algorithm for the general case that only identifies a condition describing a small fraction of the optimal condition. The former, sparse algorithm was extended to general  $\ell_p$  losses by Hainline et al. (2019). In this work, we give an algorithm that, under some mild niceness conditions,

- *only uses polynomial time and data in the dimension and number of factors,*
- *recovers a condition that covers as much of the distribution as the optimal condition, and*
- *approximates an optimal  $t$ -term  $k$ -DNF with  $\tilde{O}(t \log \log n)$  blow-up of the error.*

We also present some synthetic data experiments illustrating the capabilities of our new algorithm.

## 1.2 Related Work

Our algorithm builds on the *list-learning algorithms* due to Charikar et al. (2017). That work aimed to learn about arbitrary small subsets of the data by producing a list of parameter values containing good estimates of good parameters for any small subset. Charikar et al. could only address tasks such as mean estimation, and as they discuss, could not obtain sufficient accuracy to use their framework for linear regression. Our algorithm, like theirs, iteratively computes local estimates for the regression parameters with consideration of their neighbors, and then (re)clusters the terms using their corresponding parameters. In each iteration there is a risk to lose a small fraction of good points. Our primary innovation lies in using a fixed family of subsets (“terms” of our conditions) as the basic units of data as opposed to individual points. Since these terms are all large enough and we can obtain higher accuracy on these terms given enough data, we can show that the algorithm will not lose any terms with high probability, and it can obtain adequate estimates of the regression parameters. Specially, our improvements upon the work of Charikar et al. (2017) are:

- We modify Charikar et al’s algorithm and analysis to operate on sets instead of points, by introducing weights and modifying the definition of neighbors.
- In our improved algorithm, we strengthen the analysis to show that it does not lose good sets in any iteration, in contrast to Charikar et al’s algorithm which indeed may lose a small fraction of good points. Consequently, whereas the original algorithm may need to terminate with a relatively inaccurate estimate of the parameters, we can potentially reach any accuracy we want.
- We incorporate pre- and post-processing to convert points to our atomic sets, and to use a covering algorithm in the end to extract a  $k$ -DNF condition on the original data space.

We stress that Charikar et al. (and subsequent works such as Diakonikolas et al. (2018)) cannot obtain a linear predictor with loss that scales with the loss of the best linear predictor on the data (i.e., that goes to zero with the “noise” rate) on account of a difference in the formulation: Charikar et al. and Diakonikolas et al. consider arbitrary subsets of the data, whereas we only consider subsets described by  $k$ -DNFs. Diakonikolas et al. show that even for the simpler problem of mean estimation, one can only guarantee loss that scales polylogarithmically with the density of the set for which we estimate the mean.

Our problem is similar in spirit to work in *robust statistics* (Huber, 1981; Rousseeuw & Leroy, 1987), with the key distinction that robust statistics assumes that the outliers comprise a minority of the data. By contrast here, the vast majority of the data may be “outliers.” Another setting in

this vein is *learning with rejection* in which we think of a predictor as having the option to “abstain” from making a prediction. In most cases, the strategy for deciding when to abstain is based on some measure of “confidence” of the prediction—for example, this is how El-Yaniv and Weiner conceived of such a linear regression task (El-Yaniv & Wiener, 2012). The difference is that these methods do not generally produce a “nice” description of the region on which they will make a prediction. On the other hand, Cortes et al. (2016) considered a version of the task in which the prediction region is constrained to come from a fixed family of nice rules, like our version. The difference is that Cortes et al. do not seek to achieve given rates of coverage or error, but rather posit that abstaining from prediction has a known, fixed cost relative to the cost of an error, and seek to minimize this overall cost. In our work, we obtain a description of the regions where we will predictor or abstain as a Boolean formula like Cortes et al. (but not previous works in this area), and also give simultaneous guarantees on the overall coverage and loss unlike Cortes et al., who can only bound a weighted sum of the two. Finally, algorithms such as RANSAC (Fischler & Bolles, 1981) similarly find a dense linear relation among a subset of the points when one exists, but these algorithms scale exponentially with the dimension, and like learning with rejection, do not obtain a rule characterizing which points will satisfy the linear relationship.

There are many other works that fit the data using multiple linear rules, such as linear mixed models, segmented regression, or piecewise linear regression. These methods are similar in that is similar in that the portion of the data fit by an individual linear rule may be small. The distinction is generally that they seek to model the entire data distribution with linear rules, i.e., they cluster the data and minimize the total regression loss over all clusters. By contrast, we only seek a small fraction (say, 10%) of the data where a good linear regression fit exists. Specifically,

- *Linear mixed models* (McCulloch & Searle, 2001; Jiang, 2007) simply view the data set as a mixture of data fitting linear rules. Such work usually assumes that all or most of the data belongs to one of the mixture components. Moreover, in such models (as with learning with a reject option) the components do not come attached with a rule describing their domain. To use such models for prediction, one has to try the various linear rules to see which gives a small residual, which may not be so well behaved.
- *Generalized cluster-wise linear regression* (Park et al., 2017) describes data in terms of “entities,” where the goal is to partition these entities into clusters so that the overall total regression loss is minimized. Thus, unlike our work, these methods seek to fit the entire data set with a linear rule in some cluster. It is similar in that our terms, defined by Boolean attributes, can be viewed as entities in their setting. But, since each combination of Boolean attributes defines a term, the total number of terms could be very large (say, a million), and all these terms’ data points are defined in the same space, so they are also quite different in practice.
- In *Gaussian process classification using random decision forests* (Fröhlich et al., 2012), Fröhlich et al. propose to learn a random decision forest whose leaf nodes use a Gaussian process classifier. This algorithm also views the data as having been sampled from a linear mixture model, where again we only seek to find a small cluster. Again, this method does not produce a nice rule describing which rule to use for prediction.
- Work on *regression trees* such as by Quinlan (1992) may also be seen as finding a family of nice rules (i.e., the branches of the tree) such that on the partitions described by these rules, the data is fit by linear predictors. Again, we differ in not seeking to find linear predictors for

the entire data distribution. Furthermore, while conjunctive splits are surely nice, it seems to be intractable to identify good conjunctions, even if we are only seeking one; see Juba (2017) for details. In any case, no guarantees are known for such methods, and it is natural to conjecture that they cannot be guaranteed to work.

## 2 Preliminaries and Definitions

We suppose we have a data set consisting of  $N$  examples, where each example has three kinds of attributes: a vector of  $n$  Boolean attributes  $\mathbf{x}$ , a vector of  $d$  real-valued attributes  $\mathbf{y}$ , and a real-valued target attribute  $z$  that we wish to predict. For example, in our cancer prediction setting, we have an example  $(\mathbf{x}, \mathbf{y}, z)^{(i)}$  for each  $i$ th patient in which:

1.  $\mathbf{x}^{(i)}$  is a vector of Boolean demographic properties that describe patient  $i$  (e.g., adult, city dweller)
2.  $\mathbf{y}^{(i)}$  is vector of continuous risk factors for patient  $i$ , (e.g., rate of smoking, radiation exposure, physical activity)
3.  $z^{(i)}$  represents the variable of interest such as probability of developing a certain kind of cancer in the next ten years.

When there is no confusion, we will use  $\mathbf{x}^{(i)}$  to denote the whole point  $(\mathbf{x}, \mathbf{y}, z)^{(i)}$ .

For example,  $(y_1, y_2, y_3, z)^{(2)} = (60, 100, 18, .1)$  could represent patient 2 with height 60 inches, weight 100 lbs, age 18, and a 10% probability of developing cancer in ten years.

Note that we can obtain Boolean attributes  $\mathbf{x}$  from continuous attributes  $\mathbf{y}$  by using binning and splits, similar to decision trees. For example, we can define a new Boolean attribute such as  $x = "y \leq a?"$  for some quantile of the data  $a$ . Conversely, Boolean variables can also be viewed as regression factors.

We wish to find a linear rule  $\langle \mathbf{w}, \mathbf{y} \rangle$  to predict  $z$ . We would typically achieve this by minimizing the loss  $\|\langle \mathbf{w}, \mathbf{y} \rangle - z\|_2$  averaged over the data. But, it's common that there doesn't exist a good linear rule for the whole data set. We propose to find a subset (a condition or “cluster”  $\mathbf{c}$ ), such that there exists good fit on  $\mathbf{c}$ . Of course, if we just pick any subset that fits some linear rule well, this is unlikely to be predictive. Instead, we will seek to pick out a subset according to some simple criteria. In other words, this condition should be described by some simple rule which will use the Boolean  $\mathbf{x}$  attributes. Following the previous work Juba (2017), we will use conditions represented by  $k$ -DNF (Disjunctive normal form) formulas, which is an “or” of *terms*, where each term is an “and” of at most  $k$  attributes (which we permit to be negated). Thus:

$$\mathbf{c} = t_1 \vee \dots \vee t_s \text{ for some } s, \text{ where}$$

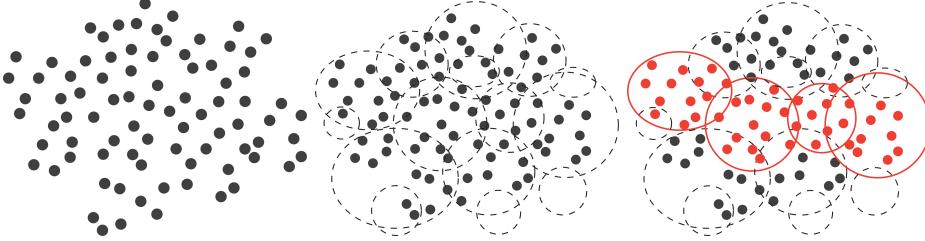
$$t_i = \ell_{i,1} \wedge \dots \wedge \ell_{i,k}, \text{ and where } \ell_{i,j} \text{ is either } x_{i,j} \text{ or } \neg x_{i,j}, \text{ for some attribute } x_{i,j}$$

and when we say  $t$ -term DNF, we mean  $s \leq t$ . We focus on  $k$ -DNF conditions since the use of other natural representations seems to result in an inherently intractable problem (Juba, 2017).

To summarize, we want to find a  $k$ -DNF condition such there exists a good linear prediction rule on the data satisfying this condition. Meanwhile, we want this condition to describe a subset of the data that isn't too small. We will demand that data belongs to it with probability at least  $\mu$ , as shown in Figure 1. Since the task is thus to choose an appropriate set of terms (defining a  $k$ -DNF), we can view the terms as  $m$  atomic sets of data.<sup>1</sup>

---

<sup>1</sup>Indeed, we will not use the structure of the terms and we could work instead with an arbitrary family of  $m$  such atomic subsets of the data, as long as we can collect enough data relative to the number of subsets.



**Figure 1:** Conditional Linear Regression

Data  $\{(x_1, \dots, x_n, y_1, z)\}$  on the  $y_1 \times z$  plane. Each term  $t$  is depicted in the  $y_1 \times z$  plane as a circle, enclosing the points satisfying that term. A  $k$ -DNF is then a union of these circles. **1:** The data space. **2:** The data and the terms. **3:** A selected subset on the data space.

**Definition 2.1 (Conditional Linear Regression)** *Conditional linear regression is the following task: given data  $\{(x_1^{(i)}, \dots, x_n^{(i)}, y_1^{(i)}, \dots, y_d^{(i)}, z^{(i)})\}_{i=1}^N$  drawn i.i.d. from a distribution  $D$ , we are to find a  $k$ -DNF condition  $\mathbf{c}$  and a regression fit  $\mathbf{w} = (w_1, \dots, w_d)$ , such that:*

- (1) *The regression loss  $\|\langle \mathbf{w}, \mathbf{y} \rangle - z\|$  is bounded for the points satisfying  $\mathbf{c}$ .*
- (2) *The condition has probability at least  $\mu$ :  $\Pr[\mathbf{c}(\mathbf{x}) = 1] \geq \mu$ .*

We will describe an algorithm that can return a pair  $(\hat{\mathbf{c}}, \hat{\mathbf{w}})$  that is close to the potential optimal solution  $(\mathbf{c}^*, \mathbf{w}^*)$ , given that the data distribution on  $\mathbf{c}$  is sufficiently nice in the following sense. First, we assume that  $\mathbf{w}^*$  gives a predictor with subgaussian residuals on  $\mathbf{c}$ . Second, we will assume that our loss function is Lipschitz. Third, we will consider the following quantity measuring how the desired conditional distribution  $D|\mathbf{c}$  varies across the terms of  $\mathbf{c}$ : if we abuse notation to let a DNF  $\mathbf{t}$  also denote a set of terms (thus, letting  $|\mathbf{t}|$  denote the number of terms in  $\mathbf{t}$ ),

$$S_{\varepsilon 0} := \max_{\mathbf{t} \subseteq \mathbf{c}: \Pr[\mathbf{t}(\mathbf{x})|\mathbf{c}] \geq \varepsilon} \frac{1}{\sqrt{|\mathbf{t}|}} \left\| \left[ \mathbb{E}[\mathbf{y}\mathbf{y}^\top | \mathbf{t}_j] - \mathbb{E}[\mathbf{y}\mathbf{y}^\top | \mathbf{c}] \right]_{\mathbf{t}_j \in \mathbf{t}} \right\|_{op}.$$

$S_{\varepsilon 0}$  measures, in a spectral sense, how different on average the distributions of sufficiently large sets of terms are from the distribution over the entire desired subset  $\mathbf{c}$ . For example, if the distribution on this subset is independent of which term of  $\mathbf{c}$  is satisfied, then  $S_{\varepsilon 0}$  is 0. Intuitively, if the distribution across these terms is very different from one another, it will be harder for our algorithm to identify that they should be fit by a common linear rule. We will need that this quantity is sufficiently small relative to the degree of (strong) convexity of the loss. More specifically, recall:

**Definition 2.2** *A function  $f : \mathcal{H} \rightarrow \mathbb{R}$  for  $\mathcal{H} \subseteq \mathbb{R}^d$  is  $\kappa$ -strongly convex if  $\forall \mathbf{w}, \mathbf{w}' \in \mathcal{H}$ ,*

$$f(\mathbf{w}') \geq \langle \mathbf{w}' - \mathbf{w}, \nabla f(\mathbf{w}) \rangle + \frac{\kappa}{2} \|\mathbf{w}' - \mathbf{w}\|_2^2.$$

Generally, for the squared error loss on a bounded space, the convexity coefficient can be viewed as a constant, when the bound  $\|\mathbf{y}\|_2 \leq B$  is fixed. Note that we will be able to enforce that our loss function is at least  $\kappa$ -strongly convex by adding a regularization term,  $\frac{\kappa}{2} \|\mathbf{w}\|_2^2$ , at the potential cost of reducing the quality of the solution. Our main theorem is of the form:

**Theorem 2.3 (Conditional  $l_2$ -Linear Regression)** *Suppose  $D$  is a joint distribution over  $x \in \{0, 1\}^n, y \in \mathcal{B} \subset \mathbb{R}^d$  and  $z \in \mathbb{R}$ . If there exists a (ideal)  $t$ -term  $k$ -DNF cluster  $\mathbf{c}^*$  and regression fit*

$\mathbf{w}^* \in \mathcal{H} \subset \mathbb{R}^{d_y}$ , where  $\mathcal{B}$  has  $l_2$  radius  $B$  and  $\mathcal{H}$  has  $l_2$  radius  $r$ , such that:

$$(1): \mathbb{E}_D[(\langle \mathbf{w}^*, \mathbf{y} \rangle - z)^2 | \mathbf{c}^*(\mathbf{x}) = 1] \leq \epsilon.$$

$$(2): \Pr[\mathbf{c}^*(\mathbf{x}) = 1] \geq \mu.$$

(3): the error  $(\langle \mathbf{w}^*, \mathbf{y} \rangle - z)$  follows a  $\sigma$ -subgaussian distribution on  $D|\mathbf{c}$

(4): the loss function  $f(\mathbf{w}, (\mathbf{y}, z)) = (\langle \mathbf{w}, \mathbf{y} \rangle - z)^2 + \frac{\kappa}{2} \|\mathbf{w}\|_2^2$  is  $L$ -lipschitz and  $\kappa$ -strongly convex,

$$\text{where } \kappa \geq \Omega\left(\frac{tS\gamma t_0 \log \frac{1}{\mu}}{\sqrt{\mu}}\right)$$

then for  $\delta, \gamma \in (0, 1)$ , using  $N = \mathcal{O}\left(\frac{B^6 d^3 \sigma^2 L^2 t^2}{\mu \gamma^4} \log(m/\delta)\right)$  examples, we can find a  $k$ -DNF  $\hat{\mathbf{c}}$  and parameters  $\hat{\mathbf{w}}$  in polynomial time, such that with probability  $(1 - \delta)$ :

$$(1): \mathbb{E}_D[(\langle \hat{\mathbf{w}}, \mathbf{y} \rangle - z)^2 | \hat{\mathbf{c}}(\mathbf{x}) = 1] \leq \mathcal{O}(t \log(\mu N)(\epsilon + \gamma))$$

$$(2): \Pr[\hat{\mathbf{c}}(\mathbf{x}) = 1] \geq (1 - \gamma)\mu.$$

We remark that the Lipschitz condition follows from the bound on  $y$  and  $w$ ,  $L \leq rB^2$ , and the number of good terms  $t$  is always at most the total number of possible terms  $m \leq n^k$ , so  $\mathcal{O}(t \log(\mu N)\epsilon) \leq \tilde{\mathcal{O}}(n^k \epsilon)$  (suppressing the other parameters).

The ideal subset  $\{x : \mathbf{c}^*(\mathbf{x}) = 1\}$  is denoted as  $I_{good}$ . When later working on the space of terms, we also use  $I_{good}$  to denote the collection of terms of the DNF  $\mathbf{c}^*$ :  $\{t_i : t_i \text{ is a term of } \mathbf{c}^*\}$ , so the number of terms in  $I_{good}$  is  $t$ . From the perspective of Charikar et al. (2017), we treat  $I_{good}$  as our “good data,” with the other points being arbitrary bad data. Our algorithm is going to suggest a list of candidate parameters  $\hat{\mathbf{w}}$ , with one of them approximating  $\mathbf{w}^*$ .

### 3 Soft Regression and Outlier Removal

Our algorithm works primarily on terms: we consider the terms to be atomic sets of data, whose weights are the number of points (probability mass) satisfying the terms. The main idea of our algorithm is to compute regression parameters  $\mathbf{w}_j$  for each term  $t_j$ , and cluster the *terms* by the distance of their *parameters*, similar to Charikar et al. (2017). Towards realizing this strategy, we need to compute approximations to the regression parameters that are not too impacted by the presence of bad data (“outliers”).

#### 3.1 Preprocessing

In this section, we show how to convert the data into a suitable form: later, we will assume the terms are disjoint and that we have an adequate number of examples to estimate the loss on each term. We will ensure these conditions by introducing duplicate points when they are shared, and by deleting terms that are satisfied by too few examples.

##### 3.1.1 Reduction to Disjoint Terms by Duplicating Points

Given  $N$  data points and  $m$  terms  $t_1, \dots, t_m$ , if we view terms as sets, our analysis will require these terms to be disjoint. A simple method is to duplicate the points for each term they are contained in. For example, if the  $i^{th}$  point  $\mathbf{x}^{(i)} = (\mathbf{x}, \mathbf{y}, z)^{(i)}$  is contained in terms  $t_a$  and  $t_b$ , then we create two points  $(\mathbf{x}, \mathbf{y}, z)^{(a,i)}$  and  $(\mathbf{x}, \mathbf{y}, z)^{(b,i)}$ , each with the same attributes  $(\mathbf{x}, \mathbf{y}, z)$  as the original point  $\mathbf{x}^{(i)}$ . After duplication, the terms are disjoint, and there will be at most  $Nm$  points. We denote the resulting number of points by  $N'$ . The size of  $I_{good}$ , changing from  $|\bigcup_{I_{good}} t_i|$  to  $\sum_{I_{good}} |t_i|$ , may also blow up with a factor ranging from 1 to  $t$ . Note that the proportion of good points  $N_{good}/N$

decreases by at most a factor of  $1/m$  since  $N'_{good}/N' \geq N_{good}/mN$ . This double counting process may skew the empirical distribution of  $I_{good}$  by up to a factor of  $t$ . Consequently, it may result in up to a  $t \leq n^k$ -factor blow-up in the error, and this is ultimately the source of the increase in loss suffered by our algorithm in Theorem 2.3. For convenience, we will use the same notation  $N$ ,  $I_{good}$  and  $\mu$  for both before and after duplication when there is no confusion.

### 3.1.2 Reduction to Adequately Sampled Terms by Deleting Small Terms

The approach of Charikar et al. (2017) can only guarantee that we obtain satisfactory estimates of the parameters for sufficiently large subsets of the data. Intuitively, this is not a significant limitation as if a term has very small size, it will not contribute much to our empirical estimates. Indeed, with high probability, the small terms (terms with size  $< \varepsilon\mu N$  for  $\varepsilon \leq \gamma/t$ ) only comprise a  $\gamma$  fraction of  $I_{good}$ . Based on this motivation, if a term has size less than  $\varepsilon\mu N$ , then we just delete it at the beginning (note here  $\varepsilon$  meaning for a fraction of data, should not be confused with  $\epsilon$  for error). Especially for a  $t$ -term DNF, not many terms could be small, so it is safe to ignore these small terms. As before, we will continue to abuse notation, using  $t$  and  $m$  for the number of terms when there is no confusion.

## 3.2 Loss Functions

In this section we define our loss functions and analyze their properties.

Given  $N$  data points and  $m$  disjoint sets (terms)  $t_1, \dots, t_m$  with size (weight)  $|t_1|, \dots, |t_m|$ , we can define a loss function for each point in the space of parameters. For each  $i$ th point, define  $f^{(i)} : \mathcal{H} \rightarrow \mathbb{R}$  by

$$f^{(i)}(\mathbf{w}) = (z^{(i)} - \langle \mathbf{w}, \mathbf{y}^{(i)} \rangle)^2$$

Similarly, we define a loss function for each of the terms  $t_j$ ,  $f_1, \dots, f_m : \mathcal{H} \rightarrow \mathbb{R}$ , as the average loss over these data points  $\{\mathbf{x}^{(i)} = (\mathbf{x}^{(i)}, \mathbf{y}^{(i)}, z^{(i)})\}$  in the term  $t_j$  (beware we abuse the notation to let  $\mathbf{x}^{(i)}$  denote the  $i$ th point).

$$\begin{aligned} f_j(\mathbf{w}) &= \frac{1}{|t_j|} \sum_{\mathbf{x}^{(i)} \in t_j} f^{(i)}(\mathbf{w}) \\ &= \frac{1}{|t_j|} \sum_{\mathbf{x}^{(i)} \in t_j} (z^{(i)} - \langle \mathbf{w}, \mathbf{y}^{(i)} \rangle)^2 \\ &\stackrel{(*)}{=} \frac{1}{|t_j|} \|\mathbf{z} - Y\mathbf{w}^\top\|_2^2 \\ &= \frac{1}{|t_j|} (\mathbf{z} - Y\mathbf{w}^\top)^\top (\mathbf{z} - Y\mathbf{w}^\top) \\ &= \frac{1}{|t_j|} \left( \mathbf{z}^\top \mathbf{z} - \mathbf{z}^\top Y\mathbf{w}^\top - \mathbf{z}^\top \mathbf{w} + \mathbf{w}^\top Y^\top \mathbf{z} \right) \\ &= \frac{1}{|t_j|} [1, \mathbf{w}] \begin{bmatrix} \mathbf{z}^\top \mathbf{z} & -\mathbf{z}^\top Y \\ -Y^\top \mathbf{z} & Y^\top Y \end{bmatrix} [1, \mathbf{w}]^\top \end{aligned}$$

Where at  $(*)$ , we write the formula in vectors and matrices. We treat  $\mathbf{z}$  as a  $|t_j| \times 1$  column vector, whose each coordinate is the  $z$  for each point in the term  $t_j$ . Similarly,  $Y$  is a  $|t_j| \times d$  matrix,

each row for a point and  $\mathbf{w}$  is  $1 \times d$  row vector. One of the advantage of our formulation is that the loss function for each term can be eventually written as  $f_j(\mathbf{w}) = [1, \mathbf{w}]A[1, \mathbf{w}]^\top$ , where  $A$  is a  $(d+1) \times (d+1)$  matrix. We can pre-compute this quadratic loss matrix  $A$  so that the running time of the main algorithm is independent of the number of data points, and is thus a function only of the number of terms and dimension for our regression problem.

Note these loss functions are stochastic, depending on the sample from the distribution  $(\mathbf{x}, \mathbf{y}, z) \sim D$ . That is, the true loss for a fixed term  $t_j$  is:

$$\mathbb{E}[f_j(\mathbf{w})] = \mathbb{E}_{\mathbf{x}^{(i)}}[(z^{(i)} - \langle \mathbf{w}, \mathbf{y}^{(i)} \rangle)^2 | t_j].$$

Similarly, for  $I_{good}$ , we define the loss function

$$f_{I_{good}}(\mathbf{w}) = \frac{1}{|I_{good}|} \sum_{\mathbf{x}^{(i)} \in I_{good}} f^{(i)}(\mathbf{w}).$$

Let  $\bar{f}$  denote the expected loss function for points averaged over  $I_{good}$ ,

$$\bar{f}(\mathbf{w}) = \mathbb{E}[f_{I_{good}}].$$

Then the optimal  $\mathbf{w}^*$  is defined as

$$\mathbf{w}^* := \arg \min_{\mathbf{w}} \bar{f}(\mathbf{w}).$$

Our ultimate goal is to find  $\hat{\mathbf{w}}$  that minimizes  $\bar{f}(\hat{\mathbf{w}})$ , but the difficulty is that  $\bar{f}$  is unknown (since  $I_{good}$  is unknown). To overcome this barrier, instead of directly minimizing  $\bar{f}(\hat{\mathbf{w}})$ , we try to find a parameter  $\hat{\mathbf{w}}$  such that  $\bar{f}(\hat{\mathbf{w}}) - \bar{f}(\mathbf{w}^*)$  is small. Once we get a close approximation  $\hat{\mathbf{w}}$ , we can use a greedy covering algorithm to find a good corresponding condition  $\hat{\mathbf{c}}$ .

In summary, we reformulate our problem in terms of these new loss functions as follows:

**Definition 3.1 (Restatement of conditional linear regression problem)** *Given  $D$  a linear regression distribution over points  $\{\mathbf{x}^{(i)} := (\mathbf{x}^{(i)}, \mathbf{y}^{(i)}, z^{(i)})\}_{(i)=1}^N$ , and  $\{t_j\}_{j=1}^m$  predefined disjoint subsets (terms), let  $I_{good}$  be the (unknown) target collection corresponding to  $\mathbf{c}^* = \bigcup_{t_j \in \mathbf{c}^*} t_j$  with probability mass  $\Pr[\mathbf{x} \in I_{good}] \geq \mu$ , and  $\bar{f}$  be the regression loss over  $I_{good}$ . If there exists a linear regression fit  $\mathbf{w}^*$  such that:*

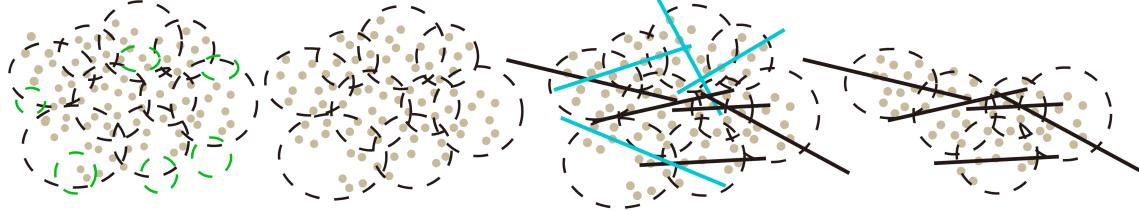
$$\mathbb{E}_D[\bar{f}(\mathbf{w}^*)] \leq \epsilon.$$

*Then we want to find a  $\hat{\mathbf{w}}$  that approximates  $\mathbf{w}^*$ :*

$$\mathbb{E}_D[\bar{f}(\hat{\mathbf{w}})] \leq \gamma + \epsilon.$$

### 3.3 Main Optimization Algorithm

The main algorithm is an alternating-minimization-style algorithm: given a soft choice of which terms are outliers, we let each term choose a local set of regression parameters that are collectively regularized by the trace of their enclosing ellipsoid. Then, given these local regression parameters, we update our scoring of outliers by examining which terms find it difficult to assemble a coalition of sufficiently many “neighboring” terms whose parameters are, on average, close to the given term. We repeat the two until we obtain a sufficiently small enclosing ellipsoid for the collection of regression parameters.



**Figure 2:** Algorithm 1

**1:** The original data space with the terms. **2:** Delete the small terms and duplicate points. **3:** Compute the best regression parameter  $\hat{\mathbf{w}}_i$  for each term. **4:** Meanwhile iteratively downweight these terms whose  $\hat{\mathbf{w}}_i$  have large error on their neighbor terms.

### 3.3.1 Semidefinite Programming for Soft Regression

Following Charikar et al. (2017), we now present Algorithm 1 for approximating the regression parameters. We assign “local” regression parameters  $\mathbf{w}_i$  for each term  $t_i$ , and use a semi-definite program (SDP) to minimize the total loss  $\sum |t_i|f_i(\mathbf{w}_i)$  with regularization to force these parameters to be close to each other. Following each iteration, we use Algorithm 2 to remove outliers, by decreasing the weight factors  $c_i$  for those terms without enough neighbors. The process is illustrated in Figure 2. Intuitively, if there exists a good linear regression fit  $\mathbf{w}^*$  on  $I_{good}$ , then for each term  $t_i \in I_{good}$ ,  $f_i(\mathbf{w}^*)$  should be small. Therefore, we can find a small ellipse  $Y$  (or  $\mathcal{E}_Y$ ) bounding all parameters for the terms in  $I_{good}$  if the center of  $Y$  is close to  $\mathbf{w}^*$ . The SDP will find such an ellipse bounding the parameters while minimizing the weighted total loss.

---

**Algorithm 1:** Soft regression algorithm

---

**Input:** terms  $t_{1:m}$

**Output:** parameters  $\hat{\mathbf{w}}_{1:m}$  and a matrix  $\hat{Y}$

Initialize  $c_{1:m} \leftarrow (1, \dots, 1)$ ,  $\lambda \leftarrow \frac{\sqrt{8\mu N t S}}{r}$

**repeat**

    Let  $\hat{\mathbf{w}}_{1:m}, \hat{Y}$  be the solution to SDP:

$$\begin{aligned} & \underset{\mathbf{w}_1, \dots, \mathbf{w}_m, Y}{\text{minimize}} \quad \sum_{i=1}^m c_i |t_i| f_i(\mathbf{w}_i) + \lambda \text{tr}(Y) \\ & \text{subject to} \quad \mathbf{w}_i \mathbf{w}_i^\top \preceq Y \text{ for all } i = 1, \dots, m. \end{aligned} \tag{1}$$

**if**  $\text{tr}(\hat{Y}) > \frac{6r^2}{\mu}$  **then**

$c \leftarrow \text{UpdateWeights}(c, \hat{\mathbf{w}}_{1:m}, \hat{Y})$

**end if**

**until**  $\text{tr}(\hat{Y}) \leq \frac{6r^2}{\mu}$

**Return**  $\hat{\mathbf{w}}_{1:m}, \hat{Y}$

---

Formally, in the SDP 1 (in Algorithm 1),  $Y$  is a  $d \times d$ -dimensional matrix (recall  $d$  is the dimension for  $\mathbf{y}$  and  $\mathbf{w}$ ). We bound the parameters  $\mathbf{w}_i$  with the ellipse  $Y$  by imposing the semidefinite constraint  $\mathbf{w}_i \mathbf{w}_i^\top \preceq Y$ , which is equivalent to letting  $\begin{bmatrix} Y & \mathbf{w}_i \\ \mathbf{w}_i^\top & 1 \end{bmatrix} \succeq 0$ , saying that  $\mathbf{w}_i$  lies within the ellipse

centered at 0 defined by  $Y$ . The regularization  $\text{tr}(Y)$  of the SDP penalizes the size of the ellipse, making the various parameter copies  $\mathbf{w}_i$  lie close to each other.

### 3.3.2 Removing Outliers

The terms not in  $I_{good}$  may have large loss for the optimal parameters  $\mathbf{w}^*$ , and therefore make the total loss in SDP 1 large. To remove these bad terms, we assign a weight factor  $c_i \in (0, 1)$  for each term  $t_i$  and down weight these terms with large loss, as shown in Algorithm 2.

---

**Algorithm 2:** Algorithm for updating outlier weights

---

**Input:**  $c, \hat{\mathbf{w}}_{1:m}, \hat{Y}$ .

**Output:**  $c'$ .

**for**  $i = 1$  **to**  $m$  **do**

    Let  $\tilde{\mathbf{w}}_i$  be the solution to

$$\begin{aligned} & \underset{\tilde{\mathbf{w}}_i, a_{i1}, \dots, a_{im}}{\text{minimize}} \quad f_i(\tilde{\mathbf{w}}_i) \\ & \text{subject to} \quad \tilde{\mathbf{w}}_i = \sum_{j=1}^m a_{ij} \hat{\mathbf{w}}_j, \quad \sum_{j=1}^m a_{ij} = 1 \\ & \quad 0 \leq a_{ij} \leq \frac{2}{\mu N} |t_j|, \quad \forall j \end{aligned} \tag{2}$$

$$z_i \leftarrow f_i(\tilde{\mathbf{w}}_i) - f_i(\hat{\mathbf{w}}_i)$$

**end for**

$$z_{max} \leftarrow \max\{z_i \mid c_i \neq 0\}$$

$$c'_i \leftarrow c_i \cdot \frac{z_{max} - z_i}{z_{max}}$$

**Return**  $c'$

---

In Algorithm 2, we solve an SDP for each term to find its best  $\mu N$  neighbor points and compute the “average” parameter  $\tilde{\mathbf{w}}_i$  over the neighborhood.  $\tilde{\mathbf{w}}_i$  is a linear combination of its neighbors’ parameters:  $\tilde{\mathbf{w}}_i = \sum_{j=1}^m a_{ij} \hat{\mathbf{w}}_j$ , minimizing the term’s loss  $f_i(\tilde{\mathbf{w}}_i)$ . Intuitively, if a term is a good term, i.e.  $t_i \in I_{good}$ , then its parameter  $\tilde{\mathbf{w}}_i$  should be close to the average of parameters of all terms in  $I_{good}$ ,  $\mathbf{w}_i \approx \sum_{I_{good}} \frac{|t_j|}{|I_{good}|} \mathbf{w}_j$ . In the SDP for  $t_i$ , we define coefficients  $a_{ij}$  to play the role of  $\frac{|t_j|}{|I_{good}|}$ . These coefficients  $\{a_{ij}\}$  are required to sum to 1, i.e.  $\sum_{j=1}^m a_{ij} = 1$ , and each should not be larger than  $\frac{|t_j|}{|I_{good}|} \sim 2 \frac{|t_j|}{\mu N}$ . At a high level, the SDP computes the best neighbors for  $t_i$  by assigning  $\{a_{ij}\}$ , so that the average parameter  $\tilde{\mathbf{w}}_i$  over the neighbors minimizes  $f_i$ . If a term is bad, it is hard to find such good neighbors, so if the loss  $f_i(\tilde{\mathbf{w}}_i)$  is much larger than the original loss, then we consider the term to be an outlier, and down-weight its weight factor  $c_i$ .

## 3.4 A Bound on the Loss That Is Linear in the Radius of Parameters

Similarly to Charikar et al. (2017), we obtain a theorem saying the algorithm will return meaningful outputs on  $I_{good}$ . The main change is that we use terms instead of points. In other words, we generalize their arguments from unit-weight points to sets with different weights. And based on a

spectral norm analysis, we show the bound will shrink linearly with the radius of the set of candidate parameters, as long as we have enough data.

First, to estimate the losses by their inputs, we introduce the gradient  $\nabla f$ . By the convexity of  $f$ , we have  $(f(\mathbf{w}) - f(\mathbf{w}^*)) \leq \langle \nabla f(\mathbf{w}), \mathbf{w} - \mathbf{w}^* \rangle$ . Note that  $\|\mathbf{w} - \mathbf{w}^*\|$  is bounded by  $2r$ , where  $r := \max \|\mathbf{w}\|_2$ . We will need to bound the gradient as well.

To bound the loss functions, we use the spectral norm of gradients:

$$S := \max_{\mathbf{w} \in \mathcal{H}} \frac{1}{\sqrt{t}} \left\| \left[ \nabla f_j(\mathbf{w}) - \nabla \bar{f}(\mathbf{w}) \right]_{j \in I_{good}} \right\|_{op}$$

where  $\|\cdot\|_{op}$  is the spectral norm (operator 2-norm) of the matrix whose rows are gradients of loss functions in  $I_{good}$ :  $(\nabla f_i(\mathbf{w}) - \nabla \bar{f}(\mathbf{w}))_{i \in I_{good}}$ .  $S$  measures the difference between the gradient of loss functions of terms in  $I_{good}$ :  $\nabla f_i(\mathbf{w})$  and gradient of average loss on  $I_{good}$ :  $\nabla \bar{f}(\mathbf{w})$ . At a high level, this bound tells us how bad these loss functions could be. We note that since the gradient is a linear operator, this quantity is invariant to regularization of the loss functions.

As shown by Charikar et al. (2017), if  $\nabla f_i - \nabla \bar{f}$  is a  $\sigma_{\nabla f}$  sub-gaussian distribution, then  $S = \mathcal{O}(\sigma_{\nabla f})$ , generally a constant. Although a constant bound  $\mathcal{O}(\sigma_{\nabla f})$  is good for their purposes – mean estimation – it is too weak for linear regression. In the sequel we will show  $S$  is going to shrink as the radius of parameters  $r$  decreases.

For linear regression,  $f_j(\mathbf{w}) := \frac{1}{|t_j|} \sum_{\mathbf{x}^{(i)} \in t_j} f^{(i)}(\mathbf{w})$ , and  $\nabla f_j(\mathbf{w}) = \frac{1}{|t_j|} \sum_{\mathbf{x}^{(i)} \in t_j} \nabla f^{(i)}(\mathbf{w})$ , where for each point  $\nabla f^{(i)}(\mathbf{w}) = 2(\mathbf{w}^\top \mathbf{y}^{(i)} - z^{(i)}) \mathbf{y}^{(i)}$ . If we assume  $z^{(i)} = \mathbf{w}^{*\top} \mathbf{y}^{(i)} + \epsilon^{(i)}$ , and the residual  $\epsilon^{(i)}$  (a subgaussian, e.g., from  $\mathcal{N}(0, \sigma_\epsilon^2)$ ) is independent of  $\mathbf{y}^{(i)}$ , then

$$\begin{aligned} \nabla f^{(i)}(\mathbf{w}) &= 2(\mathbf{w}^\top - \mathbf{w}^{*\top}) \mathbf{y}^{(i)} \mathbf{y}^{(i)\top} + \epsilon^{(i)} \mathbf{y}^{(i)} \\ \nabla f_j(\mathbf{w}) &= \frac{1}{|t_j|} \left( 2(\mathbf{w}^\top - \mathbf{w}^{*\top}) \sum_{\mathbf{x}^{(i)} \in t_j} \mathbf{y}^{(i)} \mathbf{y}^{(i)\top} + \sum_{\mathbf{x}^{(i)} \in t_j} \epsilon^{(i)} \mathbf{y}^{(i)} \right) \end{aligned}$$

Similarly, we can write the target function as:

$$\nabla \bar{f}(\mathbf{w}) = 2(\mathbf{w}^\top - \mathbf{w}^{*\top}) \mathbb{E}[\mathbf{y} \mathbf{y}^\top]$$

So the difference of the gradients is actually:

$$(\nabla f_j(\mathbf{w}) - \nabla \bar{f}(\mathbf{w})) = 2(\mathbf{w}^\top - \mathbf{w}^{*\top}) \left( \sum_{\mathbf{x}^{(i)} \in t_j} \frac{\mathbf{y}^{(i)} \mathbf{y}^{(i)\top}}{|t_j|} - \mathbb{E}[\mathbf{y} \mathbf{y}^\top] \right) + \sum_{\mathbf{x}^{(i)} \in t_j} \frac{\epsilon^{(i)} \mathbf{y}^{(i)}}{|t_j|}$$

The first term is going to shrink as  $(\mathbf{w}^\top - \mathbf{w}^{*\top})$  decreases. (If we draw enough data,  $\frac{1}{|t_j|} \sum_{t_j} \mathbf{y}^{(i)} \mathbf{y}^{(i)\top} \rightarrow \mathbb{E}[\mathbf{y}^{(i)} \mathbf{y}^{(i)\top} | t_j]$ , so we'll be able to regard the other factor as a fixed “scaling.”) The second term approaches zero as we draw more data,  $\frac{1}{|t_j|} \sum_{t_j} \epsilon^{(i)} \mathbf{y}^{(i)} \rightarrow 0$ . So given that we have drawn enough data, we will be able to bound each row of  $S$  by the radius  $r := \max_{\mathbf{w}} \|\mathbf{w}\|_2$  and similarly for the whole matrix.

More concretely, if we define  $S_0 := \left\| \left[ \frac{1}{|t_j|} \sum_{t_j} \mathbf{y}^{(i)} \mathbf{y}^{(i)\top} - \mathbb{E}[\mathbf{y} \mathbf{y}^\top] \right]_{I_{good}} \right\|_{op}$ , then we find  $S = \mathcal{O}(r S_0)$ . Note,  $S_0$  is fixed given the data, and thus remains constant across iterations. Furthermore,  $S_0$  concentrates around  $\left\| \left[ \mathbb{E}[\mathbf{y}^{(i)} \mathbf{y}^{(i)\top} | t_j] - \mathbb{E}[\mathbf{y} \mathbf{y}^\top] \right]_{I_{good}} \right\|_{op}$  and can thus be bounded. Therefore, the bound on  $S$  we can guarantee will decrease when we take more points. We know  $\frac{1}{|t_j|} \sum_{t_j} \epsilon^{(i)}$  can be

bounded with a simple sub-gaussian tail bound :  $\Pr[\frac{1}{|t_j|} \sum \epsilon \geq \tau] \leq \exp[-\frac{2\tau^2}{\sigma_\epsilon^2/|t_j|}]$ . Plugging in  $\tau \leftarrow r$ , and fixing  $\delta$ , we find that as long as the number of examples  $|t_j| \geq \sigma_\epsilon^2 \log(1/\delta)/2r^2$ , then  $\sum_{t_j} \epsilon^{(i)} \leq r$  with probability  $1 - \delta$ . Taking a union bound over  $\delta \leftarrow \delta/t$ , it suffices to take  $|t_j| \geq \sigma_\epsilon^2 \log(m/\delta)/2r^2$ , and thus  $N = \mathcal{O}(\sigma_\epsilon^2 \log(m/\delta)/\epsilon \mu r^2)$ . In summary, we obtain

**Lemma 3.2** *For  $N = \mathcal{O}(\sigma_\epsilon^2 \log(m/\delta)/\epsilon \mu r^2)$  example points, with probability  $1 - \delta$  the spectral norm of the gradients  $S$  is bounded by a linear function of the radius  $r := \max_{\mathbf{w}} \|\mathbf{w}\|_2$ , i.e.,  $S = \mathcal{O}(rS_0)$ .*

### 3.5 Analysis of Main Optimization Algorithms 1 and 2

Let  $\hat{\mathbf{w}}_{1:m}$  be the outputs from Algorithm 1. We define the weighted average parameter of terms from  $I_{good}$  as  $\hat{\mathbf{w}}_{avg} := (\sum_{i \in I_{good}} c_i |t_i| \hat{\mathbf{w}}_i) / (\sum_{i \in I_{good}} c_i |t_i|)$ . In this section, we aim to prove a bound on  $\bar{f}(\hat{\mathbf{w}}_{avg}) - \bar{f}(\mathbf{w}^*)$  by controlling the optimization error  $|f_i(\hat{\mathbf{w}}_i) - f_i(\mathbf{w}^*)|$  and the statistical error  $|\bar{f}(\hat{\mathbf{w}}_{avg}) - f_i(\hat{\mathbf{w}}_i)|$ . Then we prove Algorithm 2 will not decrease the weight of the good terms too much.

Theorem 3.3 says that Algorithm 1 can find a small ellipse  $\mathcal{E}_Y$  bounding its output, and the expected loss over  $\hat{\mathbf{w}}_{avg}$  is close to the expected loss of  $\mathbf{w}^*$ .

**Theorem 3.3 (Weighted Version of Theorem 4.1, Charikar et al. (2017))** *Let  $\hat{\mathbf{w}}_{1:m}, \hat{Y}$  be the output of Algorithm 1. Then,  $\bar{f}(\hat{\mathbf{w}}_{avg}) - \bar{f}(\mathbf{w}^*) \leq 18 \frac{tSr}{\sqrt{\mu}}$ . Furthermore,  $\hat{\mathbf{w}}_{avg} \in \mathcal{E}_{\hat{Y}}$  and  $\text{tr}(\hat{Y}) \leq \frac{6r^2}{\mu}$ .*

Lemma 3.4 is a basic inequality used multiple times in the analysis. It bounds the loss via  $S$ . Since the algorithm is using terms instead of points, we are suffering an additional factor- $t$  blow-up of the error compared to the original bound, which is carried through the lemmas in this section.

**Lemma 3.4** *For any  $\mathbf{w}$  and any  $\mathbf{w}_{1:n}$  satisfying  $\mathbf{w}_i \mathbf{w}_i^\top \preceq Y$  for all  $i$ , we have*

$$\left| \sum_{i \in I_{good}} c_i |t_i| \langle \nabla f_i(\mathbf{w}) - \nabla \bar{f}(\mathbf{w}), \mathbf{w}_i \rangle \right| \leq \mu N t S \sqrt{\text{tr}(Y)}. \quad (3)$$

**Proof of Lemma 3.4** Let  $F$  be the matrix whose  $i^{th}$  row is  $(\nabla f_i(\mathbf{w}_0) - \nabla \bar{f}(\mathbf{w}_0))$ , and let  $W$  be the matrix whose  $i^{th}$  row is  $\mathbf{w}_i$ . We consider only the rows  $i \in I_{good}$ , so the dimension of each matrix is  $t \times d$ . We have

$$\begin{aligned} \left| \sum_{i \in I_{good}} c_i |t_i| \langle \nabla f_i(\mathbf{w}_0) - \nabla \bar{f}(\mathbf{w}_0), \mathbf{w}_i \rangle \right| &= \text{tr}(F^\top \text{diag}(|t_i| c_i) W) \\ &\leq \|\text{diag}(|t_i|)\| \|\text{diag}(c_i) F\|_{op} \|W\|_* \end{aligned}$$

by Hölder's inequality. We can bound each part:

$$\begin{aligned} \|\text{diag}(t_i)\|_{op} &\leq \max_{t_i \in I_{good}} |t_i| \leq N_{good} \leq \mu N \\ \|\text{diag}(c)\|_{op} &\leq 1 \text{ since } c \in [0, 1] \\ \|F\|_{op} &\leq \sqrt{t} S, \text{ by the definition of } S \\ \|W\|_* &\leq \sqrt{t \text{tr}(Y)}, \text{ by Lemma 3.1 of Charikar et al. (2017)} \end{aligned}$$

Combining these, we see that  $\|diag(|t_i|c)F\|_{op}\|W\|_*$  is bounded by  $\mu tNS\sqrt{\text{tr}(Y)}$ . ■

Lemma 3.5 bounds the difference between  $f_i(\hat{\mathbf{w}}_i)$  and  $f_i(\mathbf{w}^*)$ , based on the optimality of our solution to SDP 1 in Algorithm 1. Its proof follows identically to Lemma 4.2 of Charikar et al. (2017).

**Lemma 3.5 (c.f. Lemma 4.2 of Charikar et al. (2017))** *The solution  $\hat{\mathbf{w}}_{1:m}$  to the SDP in Algorithm 1 satisfies:*

$$\sum_{i \in I_{good}} c_i |t_i| (f_i(\hat{\mathbf{w}}_i) - f_i(\mathbf{w}^*)) \leq \lambda \|\mathbf{w}^*\|_2^2. \quad (4)$$

Lemma 3.6 bounds the difference between  $f_i(\hat{\mathbf{w}}_{avg})$  and  $f_i(\hat{\mathbf{w}}_i)$ . Its proof likewise identically follows Lemma 4.3 of Charikar et al. (2017):

**Lemma 3.6 (c.f. Lemma 4.3 of Charikar et al. (2017))** *Let  $\hat{\mathbf{w}}_{avg} := (\sum_{i \in I_{good}} c_i |t_i| \hat{\mathbf{w}}_i) / (\sum_{i \in I_{good}} c_i |t_i|)$ . The solution  $\hat{\mathbf{w}}_{1:m}, \hat{Y}$  to Algorithm 1 satisfies*

$$\begin{aligned} \sum_{i \in I_{good}} c_i |t_i| (f_i(\hat{\mathbf{w}}_{avg}) - f_i(\hat{\mathbf{w}}_i)) &\leq \mu N t S \left( \sqrt{\text{tr}(\hat{Y})} + r \right), \\ \sum_{i \in I_{good}} c_i |t_i| (\bar{f}_i(\hat{\mathbf{w}}_{avg}) - \bar{f}(\mathbf{w}^*)) &\leq \sum_{i \in I_{good}} c_i |t_i| (f_i(\hat{\mathbf{w}}_{avg}) - f_i(\mathbf{w}^*)) + 2\mu N t r S. \end{aligned}$$

We next consider an analogue of Lemma 4.4 of Charikar et al. (2017). To deal with the different weights of terms, our Algorithm 2 considers the neighbors with their weights, and therefore uses different definitions of  $a$  and  $W$  (from those of Charikar et al. (2017)) in the analysis. Lemma 3.7 bounds  $\text{tr}(Y)$  and the difference between  $f_i(\tilde{\mathbf{w}}_i)$  and  $f_i(\hat{\mathbf{w}}_i)$ .

**Lemma 3.7** *For  $\tilde{\mathbf{w}}_i$  as obtained in Algorithm 2,  $\tilde{Y} := \frac{2}{\mu N} \hat{W} \hat{W}^\top$ , and  $W := [\sqrt{|t_1|} \mathbf{w}_1, \dots, \sqrt{|t_m|} \mathbf{w}_m]$ , we have*

$$\tilde{\mathbf{w}}_i \tilde{\mathbf{w}}_i^\top \preceq \tilde{Y}$$

for all  $i$ , and also

$$\text{tr}(\tilde{Y}) \leq \frac{2r^2}{\mu}$$

In addition:

$$\text{tr}(\tilde{Y}) \leq \frac{2r^2}{\mu} + \frac{1}{\lambda} \left( \sum_{i=1}^m c_i |t_i| (f_i(\tilde{\mathbf{w}}_i) - f_i(\hat{\mathbf{w}}_i)) \right)$$

**Proof of Lemma 3.7** Let  $\tilde{\mathbf{w}}_i = \sum_{j=1}^m a_{ij} \hat{\mathbf{w}}_j$  as defined in Algorithm 2. First, we want to show

$$\tilde{\mathbf{w}}_i \tilde{\mathbf{w}}_i^\top \preceq \tilde{Y}:$$

$$\begin{aligned}
\tilde{\mathbf{w}}_i \tilde{\mathbf{w}}_i^\top &= \left( \sum_{j=1}^n a_{ij} \hat{\mathbf{w}}_j \right) \left( \sum_{j=1}^n a_{ij} \hat{\mathbf{w}}_j \right)^\top \\
&\preceq \sum_{j=1}^n a_{ij} \hat{\mathbf{w}}_j \hat{\mathbf{w}}_j^\top \\
&\preceq \sum_{j=1}^n \frac{2}{\mu N} |t_j| \hat{\mathbf{w}}_j \hat{\mathbf{w}}_j^\top \\
&= \frac{2}{\mu N} \sum_{j=1}^n |t_j| \hat{\mathbf{w}}_j \hat{\mathbf{w}}_j^\top \\
&= \tilde{Y}
\end{aligned}$$

and

$$\begin{aligned}
\text{tr}(\tilde{Y}) &= \frac{2}{\mu N} \text{tr}(\hat{W} \hat{W}^\top) \\
&= \frac{2}{\mu N} \text{tr}(diag([|t_i|]) \| \mathbf{w} \|) \\
&\leq \frac{2}{\mu N} \sum_{i=1}^m |t_i| r^2 \\
&\leq \frac{2r^2}{\mu}.
\end{aligned}$$

For the third claim, since  $(\hat{\mathbf{w}}_{1:m}, \hat{Y})$  is the optimal solution of the SDP in Algorithm 1 and  $(\tilde{\mathbf{w}}_{1:m}, \tilde{Y})$  is a feasible solution of that, we have

$$\sum_{i=1}^m c_i |t_i| f_i(\hat{\mathbf{w}}_i) + \lambda \text{tr}(\hat{Y}) \leq \sum_{i=1}^m c_i |t_i| f_i(\tilde{\mathbf{w}}_i) + \lambda \text{tr}(\tilde{Y})$$

This gives us

$$\text{tr}(\hat{Y}) \leq \frac{2r^2}{\mu} + \frac{1}{\lambda} \left( \sum_{i=1}^m c_i |t_i| (f_i(\tilde{\mathbf{w}}_i) - f_i(\hat{\mathbf{w}}_i)) \right).$$

■

Then, analogous to Corollary 4.4 of Charikar et al. (2017), we show  $\hat{\mathbf{w}}_{avg}$  can be viewed as a feasible solution to SDP in Algorithm 2, so we can bound  $(f_i(\tilde{\mathbf{w}}_i) - f_i(\hat{\mathbf{w}}_i))$  by  $(f_i(\hat{\mathbf{w}}_{avg}) - f_i(\hat{\mathbf{w}}_i))$ .

**Corollary 3.8** *If  $\sum_{i \in I_{good}} c_i |t_i| \geq \frac{\mu N}{2}$ , then*

$$\sum_{i \in I_{good}} c_i |t_i| (f_i(\tilde{\mathbf{w}}_i) - f_i(\hat{\mathbf{w}}_i)) \leq \mu N t \left( \sqrt{\text{tr}(\hat{Y})} + r \right)$$

**Proof of Corollary 3.8:** First, we show that  $\hat{\mathbf{w}}_{avg}$  is a feasible solution for the semidefinite program for  $\tilde{\mathbf{w}}$  in Algorithm 2.

By taking  $a_{ij} = \frac{c_j|t_j|}{\sum_{j' \in I_{good}} c_{j'}|t_{j'}|}$  for  $j \in I_{good}$  and 0 otherwise, we get  $a_{ij} \leq \frac{2|t_j|}{\mu N}$  since  $\sum_{j' \in I_{good}} c_{j'}|t_{j'}| \geq \frac{\mu N}{2}$ . We see

$$\hat{\mathbf{w}}_{avg} = \frac{\sum_{j \in I_{good}} c_j|t_j|\hat{\mathbf{w}}_j}{\sum_{j' \in I_{good}} c_{j'}|t_{j'}|} = \sum_{j=1}^N a_{ij}\hat{\mathbf{w}}_j$$

Then by optimality,

$$\sum_{i \in I_{good}} c_i|t_i|(f_i(\tilde{\mathbf{w}}_i) - f_i(\hat{\mathbf{w}})) \leq \sum_{i \in I_{good}} c_i|t_i|(f_i(\hat{\mathbf{w}}_{avg}) - f_i(\hat{\mathbf{w}}))$$

which is bounded by  $\mu N t S \left( \sqrt{\text{tr}(\hat{Y})} + r \right)$  by Lemma 3.6. ■

Lemma 3.9 shows  $\sum_{I_{good}} c_i|t_i|$  is large enough. In other words, Algorithm 2 will not down weight good terms too much. Its proof follows identically to Lemma 4.5 of Charikar et al. (2017).

**Lemma 3.9 (c.f. Lemma 4.5 of Charikar et al. (2017))** Suppose that  $\frac{1}{N} \sum_{i=1}^m c_i|t_i|(f_i(\tilde{\mathbf{w}}_i) - f_i(\hat{\mathbf{w}}_i)) \geq \frac{2}{\mu N} \sum_{i \in I_{good}} c_i|t_i|(f_i(\tilde{\mathbf{w}}_i) - f_i(\hat{\mathbf{w}}_i))$ . Then, the update step in Algorithm 2 satisfies

$$\frac{1}{\mu N} \sum_{i \in I_{good}} |t_i|(c_i - c'_i) \leq \frac{1}{2N} \sum_{i=1}^m |t_i|(c_i - c'_i) \quad (5)$$

Moreover, the above supposition holds if  $\lambda = \frac{\sqrt{8\mu N t S}}{r}$  and  $\text{tr}(\hat{Y}) > \frac{6r^2}{2\mu}$ .

Finally, we prove Theorem 3.3, which bounds the difference in the empirical loss of  $\hat{\mathbf{w}}_{avg}$  and  $\mathbf{w}^*$ .

**Proof of Theorem 3.3:** First, show the weights of  $I_{good}$  will never be too small. By Lemma 3.9, the invariant  $\sum_{i \in I_{good}} c_i|t_i| \geq \frac{\mu N}{2} + \frac{\mu}{2} \sum_{i=1}^m c_i|t_i|$  holds throughout the algorithm. Therefore we get  $\sum_{i \in I_{good}} c_i|t_i| \geq \frac{\mu N}{2}$ . In particular, Algorithm 1 will terminate, since Algorithm 2 zeros out at least one outlier  $c_i$  each time, and this can happen at most  $m - t$  times before  $\sum_{i \in I_{good}} c_i|t_i|$  would drop below  $\frac{\mu N}{2}$ , which we showed impossible.

Now, let  $(\hat{\mathbf{w}}_{1:m}, \hat{Y})$  be the value returned by Algorithm 1. By Lemma 3.6 we then have

$$\begin{aligned} \sum_{i \in I_{good}} c_i|t_i|(\bar{f}(\hat{\mathbf{w}}_{avg}) - \bar{f}(\mathbf{w}^*)) &\leq \sum_{i \in I_{good}} c_i|t_i|(\bar{f}(\hat{\mathbf{w}}_{avg}) - \bar{f}(\mathbf{w}^*)) + 2\mu N t S r \\ &\leq \sum_{i \in I_{good}} c_i|t_i|(\bar{f}(\hat{\mathbf{w}}_i) - \bar{f}(\mathbf{w}^*)) + 3\mu N t S r + \sqrt{6\mu} N t S r. \end{aligned}$$

By Lemma 3.5, we have  $\sum_{i \in I_{good}} c_i|t_i|(f_i(\hat{\mathbf{w}}_i) - f_i(\mathbf{w}^*)) \leq \lambda r^2$  and, by the assumption we have  $\text{tr}(\hat{Y}) \leq \frac{6r^2}{\mu}$ . Plugging in  $\lambda = \frac{\sqrt{8\mu N t S}}{r}$ , we get

$$\begin{aligned} \sum_{i \in I_{good}} c_i|t_i|(\bar{f}(\hat{\mathbf{w}}_{avg}) - \bar{f}(\mathbf{w}^*)) &\leq \lambda r^2 + 3\mu N t S r + \sqrt{6\mu} N t S r \\ &= 3\mu N t S r + (\sqrt{6} + \sqrt{8})\sqrt{\mu} N t S r \\ &\leq 9\sqrt{\mu} N t S r \end{aligned}$$

Since  $\sum_{i \in I_{good}} c_i |t_i| \geq \frac{\mu N}{2}$ , dividing through by  $\sum_{i \in I_{good}} c_i |t_i|$  yields  $\bar{f}(\hat{\mathbf{w}}_{avg}) - \bar{f}(\mathbf{w}^*) \leq 18 \frac{tSr}{\sqrt{\mu}}$ . ■

## 4 List-regression Algorithm

We finally introduce the main algorithm to cluster the terms. Again following Charikar et al. (2017), we initially use Algorithm 1 to assign a candidate set of parameters  $\hat{\mathbf{w}}_i$  for each term. In each iteration, we use Padded Decompositions to cluster the terms by their parameters, and then reuse Algorithm 1 on each cluster. After each iteration, we can decrease the radius of the ellipse containing the parameters chosen by terms in  $I_{good}$  by half. Eventually, the algorithm will be able to constrain the parameters for all of the good terms in a very small ellipse, as illustrated in Figure 3. The algorithm will then output a list of candidate parameters, with one of them approximating  $\mathbf{w}^*$ .



**Figure 3:** Algorithm 4

**1:** Run Algorithm 1. Get a  $\hat{\mathbf{w}}_i$  for each term. **2:** Cluster the terms by their parameter  $\hat{\mathbf{w}}_i$ . **3:** Iteratively re-run Algorithm 1 on each cluster and re-cluster the terms, so that the  $\hat{\mathbf{w}}_i$  of  $I_{good}$  gradually get closer. **4:** Finally terminate by picking a “good enough” cluster.

### 4.1 Padded Decomposition

Padded Decomposition is a randomized clustering technique developed by Fakcharoenphol et al. (2003). Given points  $\{\mathbf{w}_1, \dots, \mathbf{w}_m\}$  in a metric space, a padded decomposition with parameters  $(\rho, \tau, \delta)$  is a partitioning of the points  $\mathcal{P} := \{P_i\}$  satisfying the following:

1. Each cluster  $P$  has diameter  $\rho$ ,
2. For each point  $\mathbf{w}_i$  and all  $\mathbf{w}_j$  such that  $\|\mathbf{w}_i - \mathbf{w}_j\| < \tau$ ,  $\mathbf{w}_j$  will lie in the same cluster  $P$  as  $\mathbf{w}_i$  with probability  $1 - \delta$ .

Fakcharoenphol et al. give a simple random clustering algorithm to produce padded decompositions, that uniformly samples balls with radius less than  $\rho$  from the space  $\mathcal{W} = \hat{\mathbf{w}}_{1:m}$ . Intuitively, if the radius of  $I_{good}$ ,  $\tau \ll \rho$ , then with high probability, the ball with radius  $\rho$  will contain all of  $I_{good}$ .

**Lemma 4.1 (Padded Decomposition)** *If all the elements of  $I$  have pairwise distance  $d(\hat{\mathbf{w}}_i, \hat{\mathbf{w}}_j) \leq \tau$ , and  $\rho = \frac{1}{\delta} \tau \log(\frac{1}{\mu})$ , then for the output partition  $\mathcal{P}$  of Algorithm 3, with probability least  $1 - \delta$ ,  $I$  will be contained in a single cluster  $T \in \mathbb{P}$ .*

The proof of this variant can be found in Appendix A of Charikar et al. (2017).

---

**Algorithm 3:** Padded Decomposition

---

**Input:**  $\hat{\mathbf{w}}_{1:m}, \rho, \tau$ .  
**Output:** Partition  $\mathcal{P} = \{T\}$ .  
Initialize: let  $\mathcal{P} = \emptyset$ ,  $\mathcal{W} = \hat{\mathbf{w}}_{1:m}$ . Sample  $k \sim \text{Uniform}(2, \rho)$ .  
**while**  $\mathcal{W} \neq \emptyset$  **do**  
    Sample  $i \sim \text{Uniform}(1, m)$ .  
    Let  $T \leftarrow \text{Ball}(\hat{\mathbf{w}}_i, k\tau) \cap \mathcal{W}$ .  
    Update:  $\mathcal{P} = \mathcal{P} \cup \{T\}$ .  $\mathcal{W} \leftarrow \mathcal{W} \setminus T$ .  
**end while**  
**Return:** partition  $\mathcal{P}$ .

---

## 4.2 Clustering, list regression, and conditional regression

Our algorithm uses two different kinds of clusterings of parameters to solve the list-regression problem. Given a solution to the list-regression problem, we produce a condition selecting a subset of the data (as opposed to a subset of the parameters). To reduce confusion, we now give a slightly more detailed overview of these steps, before presenting the algorithm in full detail.

In Algorithm 4, we will generate multiple padded decompositions in each iteration, to ensure that with high probability most of the padded decompositions preserve all of  $I_{good}$  in a single cluster. At the end of each iteration, we will update each term’s choice of regression parameters,  $\hat{\mathbf{w}}_i$  by aggregating the padded decompositions.

Given a target radius  $r_{final}$ , we will check if the current radius  $r < r_{final}$ . If so, the algorithm will greedily find a list of candidate parameters  $\mathbf{u}_1, \dots, \mathbf{u}_s$ , where the length of the list  $s$  is at most  $\frac{2}{\mu}$ . We can show that one of  $\mathbf{u}_1, \dots, \mathbf{u}_s$  must be close to  $\mathbf{w}^*$ . Finally, we will use a greedy covering algorithm following Juba et al. (2018) to find conditions on which the linear rules  $\mathbf{u}_1, \dots, \mathbf{u}_s$  have low loss, and return the pair  $(\hat{\mathbf{w}}, \hat{\mathbf{c}})$  with at least a  $\mu$  fraction of points and the smallest regression loss.

## 4.3 Analysis of List-regression, Algorithm 4

The analysis will require a “local” spectral norm bound that gives a tighter bound for any  $\varepsilon$  fraction of points. This analysis largely follows the same outline as Charikar et al. (2017), but differs in some key details. By the local spectral bound, in each iteration, we get good estimates of  $\mathbf{w}$  for any sufficiently large subset. A key observation is that in contrast to Charikar et al. (2017), we do not “lose” points from our clusters across iterations since our terms are all large enough that they are preserved. This enables a potentially arbitrarily-close approximation of  $\mathbf{w}^*$  given enough data.

### 4.3.1 Local Spectral Norm Bound

For  $\varepsilon < 1$ , we define a local spectral norm bound  $S_\varepsilon$  on arbitrary subsets  $T$  in  $I_{good}$ , such that  $T$  takes up at least a  $\varepsilon$  fraction of  $I_{good}$  ( $N_T \geq \varepsilon N$ ). Denote the number of points in  $T$  by  $N_T$  and the number of terms by  $m_T$ . We define

$$S_\varepsilon := \max_{w \in \mathcal{H}, T \subset I_{good}} \max_{N_T \geq \varepsilon N} \frac{1}{\sqrt{m_T}} \|[\nabla f_j(\mathbf{w}) - \nabla \bar{f}(\mathbf{w})]_{j \in T}\|_{op}$$

---

**Algorithm 4:** List-regression algorithm

---

**Input:**  $m$  terms, target radius  $r_{final}$ .  
**Output:** candidate solutions  $\{\mathbf{u}_1, \dots, \mathbf{u}_s\}$  and  $\hat{\mathbf{w}}_{1:m}$ .  
Initialize  $r^{(1)} \leftarrow r$ ,  
 $\hat{\mathbf{w}}_{1:m}^{(1)} \leftarrow$  Algorithm 1 with origin 0 and radius  $r$  (all  $i = 1, \dots, m$  are “assigned” an output).  
**for**  $\ell = 1, 2, \dots$  **do**  
     $\mathcal{W} \leftarrow \{\hat{\mathbf{w}}_i^{(\ell)} \mid \hat{\mathbf{w}}_i^{(\ell)} \text{ is assigned}\}$   
    **if**  $r^{(\ell)} < \frac{1}{2}r_{final}$  **then**  
        Greedily find a maximal set of points  $\mathbf{u}_1, \dots, \mathbf{u}_s$  s.t.  
            I:  $|B(\mathbf{u}_j; 2r_{final}) \cap \mathcal{W}| \geq (1 - \varepsilon)\mu N, \quad \forall j$ .  
            II:  $\|\mathbf{u}_j - \mathbf{u}_{j'}\|_2 > 4r_{final}, \quad \forall j \neq j'$ .  
        **Return**  $\mathcal{U} = \{\mathbf{u}_1, \dots, \mathbf{u}_s\}, \hat{\mathbf{w}}_{1:m}^{(\ell)}$ .  
    **end if**  
    **for**  $h = 1$  **to**  $112 \log(\frac{\ell(\ell+1)}{\delta})$  **do**  
         $\bar{\mathbf{w}}_{1:m}(h) \leftarrow$  unassigned  
        Let  $\mathcal{P}_h$  be a  $(\rho, 2r^{(\ell)}, \frac{7}{8})$ -padded decomposition of  $\mathcal{W}$  with  $\rho = \mathcal{O}(r^{(\ell)} \log(\frac{2}{\mu}))$ .  
        **for**  $T \in \mathcal{P}_h$  **do**  
            Let  $B(u, \rho)$  be a ball containing  $T$ . Run Algorithm 1 on  $\mathcal{H} \cap B(u, \rho)$ , with radius  $r = \rho$  and origin shifted to  $u$ .  
            for each  $\hat{\mathbf{w}}_i \in T$  assign  $\bar{\mathbf{w}}_i(h)$  as the outputs of Algorithm 1.  
        **end for**  
    **end for**  
    **end for**  
    **for**  $i = 1$  **to**  $m$  **do**  
        Find a  $h_0$  such that  $\|\bar{\mathbf{w}}_i(h_0) - \bar{\mathbf{w}}_i(h)\|_2 \leq \frac{1}{3}r^{(\ell)}$  for at least  $\frac{1}{2}$  of the  $h$ 's.  
         $\hat{\mathbf{w}}_i^{(\ell+1)} \leftarrow \bar{\mathbf{w}}_i(h_0)$  (or “unassigned” if no such  $h_0$  exists)  
    **end for**  
     $r^{(\ell+1)} \leftarrow \frac{1}{2}r^{(\ell)}$   
**end for**

---

Similar to the analysis of  $S$ ,  $S_\varepsilon = \mathcal{O}(rS_{\varepsilon 0})$ , where recall,  $S_{\varepsilon 0} := \max_{T: N_T \geq \varepsilon N} \left\| \left[ \frac{1}{|t_j|} \sum_{i \in t_j} \mathbf{y}^{(i)} \mathbf{y}^{(i)\top} - \mathbb{E}[\mathbf{y} \mathbf{y}^\top] \right]_{j \in T} \right\|_{op}$  is bounded independently of  $r^{(\ell)}$ . We denote the value of  $S_\varepsilon$  in the  $\ell^{th}$  iteration by  $S_\varepsilon^{(\ell)}$ , where  $S_\varepsilon^{(\ell)} = \mathcal{O}(r^{(\ell)} S_{\varepsilon 0})$ .

With the local spectral norm bound for the gradients, we can obtain the local version of Lemma 3.6

**Lemma 4.2 (c.f. Lemma 5.2 of Charikar et al. (2017))** *Let the weights  $b_i \in [0, 1]$  satisfy  $\sum_{i \in I_{good}} b_i |t_i| \geq \varepsilon \mu N$ , and define  $\hat{\mathbf{w}}_{avg}^b := \sum_{i \in I_{good}} b_i |t_i| \hat{\mathbf{w}}_i / \sum_{i \in I_{good}} b_i |t_i|$ . Then the output of Algorithm 1,  $\hat{\mathbf{w}}_{1:m}, \hat{Y}$  satisfies*

$$\sum_{i \in I_{good}} b_i |t_i| \left( f_i(\hat{\mathbf{w}}_{avg}^b) - f_i(\hat{\mathbf{w}}_i) \right) \leq b_i |t_i| \langle \nabla f_i(\hat{\mathbf{w}}_{avg}^b), \hat{\mathbf{w}}_{avg}^b - \hat{\mathbf{w}}_i \rangle \leq \left( \sum_{i \in I_{good}} b_i |t_i| \right) t S_\varepsilon \left( \sqrt{\text{tr}(\hat{Y})} + r \right)$$

Moreover, for any  $\mathbf{w}, \mathbf{w}' \in \mathcal{H}$ , we have

$$\left| \sum_{i \in I_{good}} b_i |t_i| (\bar{f}(\mathbf{w}) - \bar{f}(\mathbf{w}')) - \sum_{i \in I_{good}} b_i |t_i| (f_i(\mathbf{w}) - f_i(\mathbf{w}')) \right| \leq 2 \left( \sum_{i \in I_{good}} b_i |t_i| \right) t r S_\varepsilon.$$

The proof is similar to Lemma 3.6.

#### 4.3.2 Proof of the Main Theorem

We can now state and prove our main theorem for list regression. As noted at the outset, we will need to assume that the distribution over  $I_{good}$  is sufficiently similar relative to the degree of (strong) convexity of the loss.

**Theorem 4.3** *Let any  $r_{final}$  and  $\delta, \varepsilon \leq \frac{1}{2}$  be given. Suppose that the loss functions  $f_i$  are  $\kappa$ -strongly convex and  $S_{\varepsilon 0} \leq \mathcal{O}\left(\frac{\kappa\sqrt{\mu}}{t \log(1/\mu)}\right)$  for all  $i \in I_{good}$ . For  $N = \mathcal{O}(\sigma_\varepsilon^2 \log(m/\delta)/\varepsilon \mu r_{final}^2)$  example points, let  $\mathcal{U}, \hat{\mathbf{w}}_{1:m}$  be the output of Algorithm 4. Then with probability at least  $1 - \delta$ ,  $\mathcal{U}$  has size at most  $\lfloor \frac{1}{(1-\varepsilon)\mu} \rfloor$ , and  $\min_{\mathbf{u} \in \mathcal{U}} \|\mathbf{u} - \mathbf{w}^*\|_2 \leq \mathcal{O}(r_{final})$ . Moreover,  $\|\hat{\mathbf{w}}_i - \mathbf{w}^*\|_2 \leq \mathcal{O}(r_{final})$  for every term  $i \in I_{good}$ .*

Towards proving Theorem 4.3, our main step is the following bound on the quality of a single iteration of Algorithm 4:

**Theorem 4.4** *For some absolute constant  $C$ , the output  $\hat{\mathbf{w}}_{1:m}$  of Algorithm 1 during Algorithm 4 satisfies*

$$\|\hat{\mathbf{w}}_i - \mathbf{w}^*\|_2^2 \leq C \cdot \frac{r^{(\ell)} t S_\varepsilon^{(\ell)}}{\kappa \sqrt{\mu}}$$

for all terms  $i \in I_{good}$ .

The key to establishing Theorem 4.4 will be to use the bound on the statistical error from Lemma 4.2 and the strong convexity of  $f_i$ .

**Lemma 4.5** *For any  $b_i \in [0, 1]$  satisfying  $\sum_{i \in I_{good}} b_i |t_i| \geq \varepsilon \mu N$ , we have*

$$\frac{\sum_{i \in I_{good}} b_i |t_i| \|\hat{\mathbf{w}}_i - \hat{\mathbf{w}}_{avg}\|_2^2}{\sum_{i \in I_{good}} b_i |t_i|} \leq \frac{2}{\kappa} (\sqrt{\text{tr}(\hat{Y})} + r) t S_\varepsilon \quad (6)$$

**Proof of Lemma 4.5** Recall that Lemma 4.2 says that for any  $b_i \in [0, 1]$  satisfying  $\sum_{i \in I_{good}} b_i |t_i| \geq \varepsilon \mu N$ , we have

$$\sum_{i \in I_{good}} b_i |t_i| \langle \nabla f_i(\hat{\mathbf{w}}_{avg}^b), \hat{\mathbf{w}}_{avg}^b - \hat{\mathbf{w}}_i \rangle \leq \left( \sum_{i \in I_{good}} b_i |t_i| \right) t S_\varepsilon \left( \sqrt{\text{tr}(\hat{Y})} + r \right)$$

By strong convexity of  $f_i$ , we have

$$\begin{aligned} 0 &\leq \sum_{i \in I_{good}} b_i |t_i| (f_i(\hat{\mathbf{w}}_{avg}^b) - f_i(\hat{\mathbf{w}}_i)) \\ &\leq \sum_{i \in I_{good}} b_i |t_i| \left( \langle \nabla f_i(\hat{\mathbf{w}}_{avg}^b), \hat{\mathbf{w}}_{avg}^b - \hat{\mathbf{w}}_i \rangle - \frac{\kappa}{2} \|\hat{\mathbf{w}}_i - \hat{\mathbf{w}}_{avg}^b\|_2^2 \right) \\ &\leq \left( \sum_{i \in I_{good}} b_i |t_i| \right) t S_\varepsilon \left( \sqrt{\text{tr}(\hat{Y})} + r \right) - \frac{\kappa}{2} \sum_{i \in I_{good}} b_i |t_i| \|\hat{\mathbf{w}}_i - \hat{\mathbf{w}}_{avg}^b\|_2^2. \end{aligned}$$

■

By applying Lemma 4.5 to  $b'_i = \frac{1}{2} \left( b_i + \frac{\sum_j b_j |t_j|}{\sum_j c_j |t_j|} c_i \right)$ , we obtain Lemma 4.6, which gives bounds in terms of  $\hat{\mathbf{w}}_{avg}$  rather than  $\hat{\mathbf{w}}_{avg}^b$ :

**Lemma 4.6** *For any  $b_i \in [0, 1]$  satisfying  $\varepsilon \mu N \leq \sum_{i \in I_{good}} b_i |t_i| \leq \sum_{i \in I_{good}} c_i |t_i|$ , we have*

$$\frac{\sum_{i \in I_{good}} b_i |t_i| \|\hat{\mathbf{w}}_i - \hat{\mathbf{w}}_{avg}\|_2^2}{\sum_{i \in I_{good}} b_i |t_i|} \leq \frac{10}{\kappa} (\sqrt{\text{tr}(\hat{Y})} + r) t S_\varepsilon \quad (7)$$

**Proof of Lemma 4.6** For convenience, let us define:  $B = \sum_{i \in I_{good}} b_i |t_i|$ ,  $C = \sum_{i \in I_{good}} c_i |t_i|$ ,  $b'_i = \frac{1}{2} \left( b_i + \frac{\sum_j b_j |t_j|}{\sum_j c_j |t_j|} c_i \right) = \frac{1}{2} b_i + \frac{1}{2} \frac{B}{C} c_i$ . Notice that  $\sum_{I_{good}} b'_i |t_i| = B$  and  $\hat{\mathbf{w}}_{avg}^{b'} = \frac{1}{2} \hat{\mathbf{w}}_{avg}^b + \frac{1}{2} \hat{\mathbf{w}}_{avg}$ .

We invoke Lemma 4.5 twice, on  $b$  and  $b'$  respectively:

$$\begin{aligned} \frac{1}{B} \sum_{i \in I_{good}} b'_i |t_i| \|\hat{\mathbf{w}}_i - \frac{1}{2} \hat{\mathbf{w}}_{avg}^b - \frac{1}{2} \hat{\mathbf{w}}_{avg}\| &= \frac{1}{B} \sum_{i \in I_{good}} b'_i |t_i| \left\| \frac{1}{2} \hat{\mathbf{w}}_i - \frac{1}{2} \hat{\mathbf{w}}_{avg}^{b'} \right\| \leq \frac{2}{\kappa} (\sqrt{\text{tr}(\hat{Y})} + r) t S_\varepsilon \\ \frac{1}{B} \sum_{i \in I_{good}} b_i |t_i| \left\| \frac{1}{2} \hat{\mathbf{w}}_i - \frac{1}{2} \hat{\mathbf{w}}_{avg}^b \right\| &\leq \frac{1}{\kappa} (\sqrt{\text{tr}(\hat{Y})} + r) t S_\varepsilon. \end{aligned}$$

Since for any  $i$ ,  $b'_i \leq \frac{1}{2} b_i$

$$\frac{1}{B} \sum_{i \in I_{good}} b_i |t_i| \|\hat{\mathbf{w}}_i - \frac{1}{2} \hat{\mathbf{w}}_{avg}^b - \frac{1}{2} \hat{\mathbf{w}}_{avg}\| \leq \frac{2}{B} \sum_{i \in I_{good}} b'_i |t_i| \|\hat{\mathbf{w}}_i - \frac{1}{2} \hat{\mathbf{w}}_{avg}^b - \frac{1}{2} \hat{\mathbf{w}}_{avg}\| \leq \frac{4}{\kappa} (\sqrt{\text{tr}(\hat{Y})} + r) t S_\varepsilon$$

Combining the two inequalities

$$\begin{aligned} \frac{1}{B} \sum_{i \in I_{good}} b'_i |t_i| \|\hat{\mathbf{w}}_i - \hat{\mathbf{w}}_{avg}^b\| &\leq \frac{1}{B} \sum_{i \in I_{good}} b'_i |t_i| \left( \|\hat{\mathbf{w}}_i - \frac{1}{2} \hat{\mathbf{w}}_{avg}^b - \frac{1}{2} \hat{\mathbf{w}}_{avg}\| + \left\| -\frac{1}{2} \hat{\mathbf{w}}_i + \frac{1}{2} \hat{\mathbf{w}}_{avg}^b \right\| \right) \\ &\leq 2 \left( \frac{4}{\kappa} (\sqrt{\text{tr}(\hat{Y})} + r) t S_\varepsilon + \frac{1}{\kappa} (\sqrt{\text{tr}(\hat{Y})} + r) t S_\varepsilon \right) \\ &= \frac{10}{\kappa} (\sqrt{\text{tr}(\hat{Y})} + r) t S_\varepsilon \end{aligned}$$

■

**Corollary 4.7** *In particular, no set of terms comprising more than a  $\varepsilon \mu$  fraction of the data (or probability weight) can have  $\|\hat{\mathbf{w}}_i - \hat{\mathbf{w}}_{avg}\|_2^2 > \frac{10}{\kappa} (\sqrt{\text{tr}(\hat{Y})} + r) t S_\varepsilon$ .*

**Proof** Consider terms  $t_1, \dots, t_q$  with  $\Pr[t_1 \vee \dots \vee t_q] > \varepsilon \mu$ . Assume for contradiction that for all of these terms,  $\|\hat{\mathbf{w}}_i - \hat{\mathbf{w}}_{avg}\|_2^2 > \frac{10}{\kappa} (\sqrt{\text{tr}(\hat{Y})} + r) t S_\varepsilon$ . We can assign  $b_i$  for each  $t_i$  such that  $\sum b_i = \varepsilon \mu$ . Then

$$\begin{aligned} \frac{\sum_{i \in I_{good}} b_i |t_i| \|\hat{\mathbf{w}}_i - \hat{\mathbf{w}}_{avg}\|_2^2}{\sum_{i \in I_{good}} b_i |t_i|} &> \frac{\sum_{i \in I_{good}} b_i |t_i| \frac{10}{\kappa} (\sqrt{\text{tr}(\hat{Y})} + r) t S_\varepsilon}{\sum_{i \in I_{good}} b_i |t_i|} \\ &= \frac{10}{\kappa} (\sqrt{\text{tr}(\hat{Y})} + r) t S_\varepsilon \end{aligned}$$

which contradicts Lemma 4.6. ■

**Key Observation** As we deleted all terms of size smaller than  $\varepsilon\mu N$ , all the remaining terms have at least  $\varepsilon\mu$  probability-weight (or  $\varepsilon\mu N$  empirical size). Therefore, every term will satisfy

$$\|\hat{\mathbf{w}}_i - \hat{\mathbf{w}}_{avg}\|_2^2 \leq \frac{10}{\kappa}(\sqrt{\text{tr}(\hat{Y})} + r)tS_\varepsilon.$$

We can subsequently obtain Theorem 4.4 by thus invoking Corollary 4.7:

**Proof of Theorem 4.4:** By Corollary 4.7, we have  $\|\hat{\mathbf{w}}_i - \hat{\mathbf{w}}_{avg}\|_2^2 \leq \frac{10}{\kappa}(\sqrt{\text{tr}(\hat{Y})} + r)tS_\varepsilon$  for all  $i \in I_{good}$ .  $\text{tr}(\hat{Y}) \leq \mathcal{O}(\frac{r^2}{\mu})$ , so  $\|\hat{\mathbf{w}}_i - \hat{\mathbf{w}}_{avg}\|_2^2 \leq \frac{10}{\kappa}(\frac{r}{\sqrt{\mu}} + r)tS_\varepsilon = \mathcal{O}(\frac{rtS_\varepsilon}{\kappa\sqrt{\mu}})$ . In addition, by Theorem 3.3,  $\bar{f}(\hat{\mathbf{w}}_{avg}) - \bar{f}(\mathbf{w}^*) \leq \mathcal{O}(\frac{rtS_\varepsilon}{\sqrt{\mu}})$ . By the strong convexity of  $\bar{f}$ ,  $\|\hat{\mathbf{w}}_{avg} - \mathbf{w}^*\|_2^2 \leq \mathcal{O}(\frac{rtS_\varepsilon}{\sqrt{\mu}})$ . We combine the bounds to obtain  $\|\hat{\mathbf{w}}_i - \mathbf{w}^*\|_2^2 \leq 2(\|\hat{\mathbf{w}}_i - \hat{\mathbf{w}}_{avg}\|_2^2 + \|\hat{\mathbf{w}}_{avg} - \mathbf{w}^*\|_2^2) \leq \mathcal{O}(\frac{rtS_\varepsilon}{\sqrt{\mu}})$  ■

Finally, using Theorem 4.4, we show the radius  $r^{(\ell)}$  (used in the  $\ell^{th}$  iteration) can be decreased by half at each iteration.

**Lemma 4.8** *In Algorithm 4, denote the set of parameters of good points of  $\ell^{th}$  iteration by  $I_{good}^{(\ell)} := \{\hat{\mathbf{w}}_i^{(\ell)} : i \in I_{good}\}$ . If  $\|\hat{\mathbf{w}}_i^{(\ell)} - \mathbf{w}^*\|_2 \leq r^{(\ell)}$  and  $S_{\varepsilon 0} \leq C' \cdot \frac{\kappa\sqrt{\mu}}{t\log(2/\mu)}$  for some constant  $C'$ , then with probability  $(1 - \frac{\delta}{\ell(\ell+1)})$  over the randomly chosen padded decompositions,  $\|\hat{\mathbf{w}}_i^{(\ell+1)} - \mathbf{w}^*\|_2 \leq \frac{1}{2}r^{(\ell)}$ .*

**Proof** We call a padded decomposition partition  $\mathcal{P}_h$  *good* if all of the terms of  $I_{good}^{(\ell)}$  lie in a single cluster of  $\mathcal{P}_h$ . Denote the set of padded decompositions where  $\mathcal{P}_h$  is good by  $H$ .

In the algorithm, we draw  $q = 112\log\frac{\ell(\ell+1)}{\delta}$  random padded decompositions with parameters  $(\rho, 2r^{(\ell)}, \frac{1}{8})$ , where  $\rho = \mathcal{O}(r^{(\ell)}\log\frac{2}{\mu})$ , so that (i) each cluster  $P$  of  $\mathcal{P}_h$  has diameter at most  $\mathcal{O}(r^{(\ell)}\log\frac{2}{\mu})$ , and (ii) for each padded decomposition and a parameter vector, all other parameter vectors within  $2r^{(\ell)}$  will lie in the same cluster with probability  $7/8$ .

Since we assume  $\|\hat{\mathbf{w}}_i^{(\ell)} - \mathbf{w}^*\|_2 \leq r^{(\ell)}$ , with probability  $\frac{7}{8}$ , all of  $I_{good}^{(\ell)}$  will lie in a single cluster, i.e. this padded decomposition  $\mathcal{P}_h$  is good. Then, by a Chernoff Bound, the total number of good padded decompositions will be larger than  $\frac{3}{4}q$  with probability  $1 - \frac{\delta}{\ell(\ell+1)}$ .

For a good padded decomposition ( $P$  is the cluster containing the terms of  $I_{good}^{(\ell)}$ ),  $\mathbf{w}^*$  is within distance  $r^{(\ell)}$  of  $P$ . Therefore, if  $P \subset B(u, \rho)$ , then  $\mathbf{w}^* \in B(u, \rho + r^{(\ell)})$ . As we run Algorithm 1 on  $B(u, \rho + r^{(\ell)})$ , Theorem 4.4 will give us:

$$\begin{aligned} \|\bar{\mathbf{w}}_i(h) - \mathbf{w}^*\|_2^2 &\leq C \cdot \frac{trS_\varepsilon}{\kappa\sqrt{\mu}} \\ \|\bar{\mathbf{w}}_i(h) - \mathbf{w}^*\|_2 &\leq \mathcal{O}\left(\sqrt{\frac{t(\rho + r^{(\ell)})S_\varepsilon^{(\ell)}}{\kappa\sqrt{\mu}}}\right) \\ &= \mathcal{O}\left(\sqrt{\frac{tr^{(\ell)}\log\frac{2}{\mu}S_\varepsilon^{(\ell)}}{\kappa\sqrt{\mu}}}\right) \end{aligned}$$

where  $\bar{\mathbf{w}}_i(h)$  is the output  $\hat{\mathbf{w}}_i$  of Algorithm 1.

We want to show  $\|\bar{\mathbf{w}}_i(h) - \mathbf{w}^*\|_2 \leq \frac{1}{6}r^{(\ell)}$ . Recall that  $S_\varepsilon^{(\ell)} \leq \mathcal{O}(S_{\varepsilon 0}r^{(\ell)})$ , so it suffices to have

$$r^{(\ell)} \cdot \sqrt{\frac{t \log \frac{2}{\mu} S_{\varepsilon 0}}{\kappa \sqrt{\mu}}} \leq \mathcal{O}(r^{(\ell)}),$$

i.e.,  $S_{\varepsilon 0} \leq \mathcal{O}\left(\frac{\kappa \sqrt{\mu}}{t \log \frac{2}{\mu}}\right)$ , which is true by hypothesis (for some suitable  $C'$ ).

Since  $\|\bar{\mathbf{w}}_i(h) - \mathbf{w}^*\|_2 \leq \frac{1}{6}r^{(\ell)}$  for any two good iterations  $h$  and  $h'$ , and each iteration is only bad with probability  $\leq \frac{1}{4}$ , we can pick  $h_0$  such that  $\|\bar{\mathbf{w}}_i(h_0) - \bar{\mathbf{w}}_i(h)\|_2 \leq \frac{1}{3}r^{(\ell)}$  is true for half of the iterations. For any good  $h$ ,

$$\begin{aligned} \|\bar{\mathbf{w}}_i(h_0) - \mathbf{w}^*\|_2 &\leq \|\bar{\mathbf{w}}_i(h_0) - \bar{\mathbf{w}}_i(h)\|_2 + \|\bar{\mathbf{w}}_i(h) - \mathbf{w}^*\|_2 \\ &\leq \frac{1}{3}r^{(\ell)} + \frac{1}{6}r^{(\ell)} \\ &\leq \frac{1}{2}r^{(\ell)} \end{aligned}$$

That is,  $\hat{\mathbf{w}}_i^{(\ell+1)} := \bar{\mathbf{w}}_i(h_0)$  is within  $\frac{1}{2}r^{(\ell)}$  of  $\mathbf{w}^*$ . ■

**Proof of Theorem 4.3:** Lemma 4.8 shows that on the  $\ell$ th iteration,  $\mathbf{w}^*$  lies in a ball of radius  $r^{(\ell)}$  around each  $\hat{\mathbf{w}}_i^{(\ell)}$  for  $i \in I_{good}$ , where  $r^{(\ell)}$  decreases by half on each iteration. In the final iteration, when  $r^{(\ell)}$  reaches the target accuracy radius  $r_{final}$ , the algorithm greedily finds disjoint balls  $B(\mathbf{u}_j; 2r_{final})$  on the parameter space  $\mathcal{W}$ , such that the corresponding terms for the covered parameters contain at least  $(1 - \varepsilon)\mu N$  points, i.e for each ball  $B$ ,

$$\sum_{\hat{\mathbf{w}}_i \in B} |t_i| \leq (1 - \varepsilon)\mu N$$

Since for all  $i \in I_{good}$ ,  $\|\hat{\mathbf{w}} - \mathbf{w}^*\| < r_{final}$ , we now argue that all of the terms in  $I_{good}$  will lie in one ball. Indeed, if no term of  $I_{good}$  is contained in any ball, then any term of  $I_{good}$  gives a candidate for the greedy algorithm to add to the list. Therefore, at least one of the good terms  $\hat{\mathbf{w}}_i, i \in I_{good}$  must be contained in some ball. Then for this ball  $B(\mathbf{u}, r_{final})$ ,

$$\begin{aligned} \|\mathbf{u} - \mathbf{w}^*\| &\leq \|\mathbf{u} - \hat{\mathbf{w}}_i\| + \|\hat{\mathbf{w}}_i - \mathbf{w}^*\| \\ &\leq 2r_{final} + r_{final} \\ &= \mathcal{O}(r_{final}) \end{aligned}$$

Since each ball contains at least  $(1 - \varepsilon)\mu$  points, there can be at most  $\lfloor \frac{1}{(1-\varepsilon)\mu} \rfloor$  such balls, which completes the proof. ■

## 5 Obtaining a $k$ -DNF Condition

Once we get outputs  $\{\mathbf{u}_1, \dots, \mathbf{u}_s\}$  from Algorithm 4, we switch from the parameter space  $\{\mathbf{w}\}$  back to the Boolean data space  $\{\mathbf{x}\}$ , to search for corresponding conditions  $\mathbf{c}$  for each candidate parameter  $\mathbf{u}_i$ . If we find a pair  $(\mathbf{u}, \mathbf{c})$  such that  $\mathbf{c}$  contains enough points and the loss  $f_{\mathbf{c}}(\mathbf{u})$  is small, we return this pair as the final solution.

Suppose  $\mathbf{u}$  is one of the candidates such that  $\|\mathbf{u} - \mathbf{w}^*\| < \mathcal{O}(r_{final}) =: \gamma$ , then  $|\bar{f}(\mathbf{u}) - \bar{f}(\mathbf{w}^*)| \leq \gamma L = \mathcal{O}(\gamma)$ , for some Lipschitz constant  $L$  (since  $f$  is just a regression loss on a bounded space, it is Lipschitz continuous). Recalling  $\bar{f}$  is nonnegative, if  $\bar{f}(\mathbf{w}^*) \leq \epsilon$ , then  $\bar{f}(\mathbf{u}) \leq \epsilon + \mathcal{O}(\gamma)$ . The above also holds for the empirical approximation to  $\bar{f}$ ,  $f_{I_{good}}$ .

## 5.1 Bounding the Double-Counting Effect

We now address the effect of our double-counting of points. Recall that we introduced a copy of a point for each term it satisfied. Observe that on  $I_{good}$ , which contains  $t$  terms, this is at most  $t$  copies. We thus obtain

**Lemma 5.1** *Let  $\mathbf{u}$  be such that  $\|\mathbf{u} - \mathbf{w}^*\| < \gamma$ . Then  $f_{I_{good}}(\mathbf{u}) \leq \frac{1}{|I_{good}|} \sum_{t_j \in \mathbf{c}^*} |t_j| f_j(\mathbf{u}) + \mathcal{O}(t\gamma) \leq t\epsilon + \mathcal{O}(t\gamma)$  (where  $f_{I_{good}}$  refers to the true empirical loss, without duplicated points).*

**Proof** Assume  $\mathbf{w}_{true}^*$  is the true optimal linear fit: ignoring the common  $1/|I_{good}|$  scaling,

$$\mathbf{w}_{true}^* := \operatorname{argmin}_{\mathbf{w}} f_{I_{good}}(\mathbf{w}) = \operatorname{argmin}_{\mathbf{w}} \sum_{i \in I_{good}} f^{(i)}(\mathbf{w})$$

and  $\mathbf{w}^*$  is the optimal linear fit for the double counted data: letting  $a^{(i)} \in [1, t]$  denote the number of copies of each point after duplication and again ignoring the  $\sum_{i \in I_{good}} a^{(i)}$  scaling factor,

$$\mathbf{w}^* := \operatorname{argmin}_{\mathbf{w}} \sum_{i \in I_{good}} a^{(i)} f^{(i)}(\mathbf{w}).$$

Now, for any  $\mathbf{w}$ , observe that  $\sum_{i \in I_{good}} f^{(i)}(\mathbf{w}) \leq \sum_{i \in I_{good}} a^{(i)} f^{(i)}(\mathbf{w})$  and

$$\sum_{i \in I_{good}} a^{(i)} f^{(i)}(\mathbf{w}) \leq (\max_{i \in I_{good}} a^{(i)}) \sum_{i \in I_{good}} f^{(i)}(\mathbf{w}) \leq t \sum_{i \in I_{good}} f^{(i)}(\mathbf{w}).$$

Therefore,

$$|I_{good}| f_{I_{good}}(\mathbf{w}^*) \leq \sum_{i \in I_{good}} a^{(i)} f^{(i)}(\mathbf{w}^*) \leq \sum_{i \in I_{good}} a^{(i)} f^{(i)}(\mathbf{w}_{true}^*) \leq t |I_{good}| f_{I_{good}}(\mathbf{w}_{true}^*)$$

where we note that for any  $\mathbf{w}$ ,  $\sum_{i \in I_{good}} a^{(i)} f^{(i)}(\mathbf{w}) = \sum_{t_j \in \mathbf{c}^*} |t_j| f_j(\mathbf{w})$ . Therefore, if  $f_{I_{good}}(\mathbf{w}_{true}^*) \leq \epsilon + \mathcal{O}(\gamma)$ , then

$$f_{I_{good}}(\mathbf{w}^*) \leq \frac{1}{|I_{good}|} \sum_{t_j \in \mathbf{c}^*} |t_j| f_j(\mathbf{w}^*) \leq t\epsilon + \mathcal{O}(t\gamma)$$

where  $|f_{I_{good}}(\mathbf{w}^*) - f_{I_{good}}(\mathbf{u})| \leq \mathcal{O}(\gamma)$  and  $|\sum_{t_j \in \mathbf{c}^*} \frac{|t_j|}{|I_{good}|} f_j(\mathbf{w}^*) - \sum_{t_j \in \mathbf{c}^*} \frac{|t_j|}{|I_{good}|} f_j(\mathbf{u})| \leq \mathcal{O}(t\gamma)$ . ■

## 5.2 Greedy Set-Cover

We have obtained a parameter vector  $\mathbf{u}$  such that the loss for each term  $f_i(\mathbf{u})$  is close to  $f_i(\mathbf{w}^*)$ . We can now use a greedy weighted partial set-cover algorithm for the *ratio objective* (optimizing the ratio of the cost to the number of elements covered), e.g., as presented and analyzed by Zhang et al. (2017) following Slavík (1997), to find the corresponding conditions  $\mathbf{c}$ . At a high level, given regression parameters  $\mathbf{u}$ , we compute the loss  $f(\mathbf{u})$  for each point, and then use the covering algorithm to find a collection of terms that cover enough points while minimizing the loss.

Specifically, we use Algorithm 5, associating the term  $t_j$  with the set  $T_j = \{\mathbf{x}^{(i)} : t_j(\mathbf{x}^{(i)}) = 1\}$  and with cost  $\omega_j = |t_j| f_j(\mathbf{u})$ . In other words, it associates with  $t_j$  the set of examples such that

---

**Algorithm 5:** Partial Greedy Algorithm

---

**Input:** finite set  $\mathcal{T} = \{T_1, \dots, T_m\}$ , costs  $\{\omega_1, \dots, \omega_m\}$ ,  $\mu, \gamma \in (0, 1]$

**Output:** Condition  $\hat{\mathbf{c}}$

- 1: Initialize  $\hat{\mathbf{c}} = \emptyset$
  - 2: **while**  $(1 - (2/3)\gamma)\mu N > \left| \bigcup_{T_j \in \hat{\mathbf{c}}} T_j \right|$  **do**
  - 3:     Choose the first  $T_j \in \mathcal{T} \setminus \hat{\mathbf{c}}$  covering at least  $\frac{\mu\gamma}{3t}N$  additional examples, that minimizes  $\omega_j/|T_j|$ , for  $T_j \in \mathcal{T} \setminus \hat{\mathbf{c}}$ .
  - 4:     Add  $T_j$  to  $\hat{\mathbf{c}}$ , set each other  $T_{j'} = T_{j'} \setminus T_j$
  - 5: **end while**
  - 6: Return  $\hat{\mathbf{c}}$
- 

$t_j(\mathbf{x}^{(i)}) = 1$ , and assigns each set the cost  $\sum_{\mathbf{x}^{(i)}, t_j(\mathbf{x}^{(i)})=1} f^{(i)}(\mathbf{u})$ . It then follows the standard greedy algorithm for weighted partial set cover on this instance, modified slightly to ignore sets that cover too few examples. This latter condition will allow us to control the size of the formula  $\hat{\mathbf{c}}$  we find as a function of  $t$  and  $\gamma$ . Following Juba et al. (2018), we will be able to leverage this bound on the size of  $\hat{\mathbf{c}}$  to obtain a better approximation ratio for the loss.

**Lemma 5.2** *Given a set of  $N = \mathcal{O}\left(\frac{1}{\mu\gamma^2}\left(\frac{t}{\gamma} + \frac{\sigma^2 L^2}{\epsilon}\right) \log \frac{m}{\delta}\right)$  points  $\{\mathbf{x}^{(i)}\}_{i=1}^N$  with weights  $f^{(i)}(\mathbf{u})$  and terms  $\{t_j\}_{j=1}^m$ , if there exists a  $t$ -term  $k$ -DNF  $\mathbf{c}^*$  that is satisfied by a  $\mu$ -fraction of the points with total loss  $\epsilon$ , then the weighted greedy set cover algorithm (Algorithm 5) finds a  $3t/\gamma$ -term  $k$ -DNF  $\hat{\mathbf{c}}$ , that is satisfied by a  $(1 - \gamma)\mu$ -fraction of the points with total loss  $\mathcal{O}(t \log(\mu N) \epsilon)$*

**Proof** Observe that since the loss is non-negative, for any  $\mathbf{w}$  the terms  $t_j$  in a formula  $\mathbf{c}$  must satisfy

$$\mathbb{E}[f(\mathbf{w})|\mathbf{c}] \Pr[\mathbf{c}(\mathbf{x}) = 1] \leq \sum_{t_j \in \mathbf{c}} \mathbb{E}[f(\mathbf{w})|t_j] \Pr[t_j(\mathbf{x}) = 1]$$

as the RHS simply counts the contribution of some points multiple times, depending on how many terms of  $\mathbf{c}$  it satisfies. The latter quantities are approximated by  $\frac{|t_j|}{N} f_j(\mathbf{w})$ . By taking  $N$  sufficiently large, we can ensure that with high probability,  $\frac{1}{N} \sum_{t_j \in \mathbf{c}^*} |t_j| f_j(\mathbf{u})$  is at most  $t\mu\epsilon + \mathcal{O}(t\mu\gamma)$  by Lemma 5.1. We note that the cost of any cover  $\mathbf{c}$  is  $\sum_{t_j \in \mathbf{c}} |t_j| f_j(\mathbf{u})$ , which is the (double-counted) loss of  $\mathbf{c}$  (up to rescaling).

It follows from an analysis by Haussler (1988) that since there are at most  $\sim m^{3t/\gamma}$  formulas consisting of at most  $\frac{3t}{\gamma}$  terms,  $\mathcal{O}\left(\frac{t}{\mu\gamma^3} \log \frac{m}{\delta}\right)$  examples suffice to guarantee that any  $\frac{3t}{\gamma}$ -term formula (out of the  $m$  possible terms) that empirically satisfies at least  $(1 - (2/3)\gamma)\mu N$  examples must be satisfied with probability at least  $(1 - \gamma)\mu$  overall, and conversely, any formula (in particular,  $\mathbf{c}^*$ ) that is true with probability at least  $\mu$  will empirically satisfy at least  $(1 - \gamma/3)\mu N$  examples. Moreover, we can obtain more generally that  $\mathcal{O}\left(\frac{\sigma^2 L^2}{\mu\epsilon\gamma^2} \log \frac{m}{\delta}\right)$  examples suffice to guarantee that the loss on all terms is estimated to within a  $(1 \pm \gamma)$ -factor. The above will simultaneously hold with probability  $1 - \delta$  for suitable constants.

Zhang et al. (2017) showed that the greedy algorithm obtains a  $3H((1 - (2/3)\gamma)\mu N)$ -approximation to the minimum weight set cover under the ratio objective, where  $H(\ell)$  denotes the  $\ell$ th harmonic number, which is  $\leq \log(\mu N) + 1$ . We have modified the algorithm slightly, to ignore sets that cover fewer than  $\frac{\mu\gamma}{3t}N$  points. Observe that if all  $t$  terms of  $\mathbf{c}^*$  fail this condition, then in total there

must be at most  $t \cdot \frac{\mu\gamma}{3t} N = \frac{\mu\gamma}{3} N$  points out of the  $(1 - \gamma/3)\mu N$  points of  $\mathbf{c}^*$  uncovered, so at least  $(1 - (2/3)\gamma)\mu N$  points are already covered by  $\hat{\mathbf{c}}$ , and thus the algorithm would already terminate. We can thus still compare the ratio of the set chosen by the greedy algorithm to the term of  $\mathbf{c}^*$  that still picks up sufficiently many additional points, with the smallest ratio. Thus, the rest of the argument remains the same and the greedy algorithm finds a formula  $\hat{\mathbf{c}}$  with

$$\frac{\sum_{t_j \in \hat{\mathbf{c}}} |t_j| f_j(\mathbf{u})}{|\{\mathbf{x}^{(i)} : \hat{\mathbf{c}}(\mathbf{x}^{(i)}) = 1\}|} \leq 3H(\mu N) \frac{\sum_{t_j \in \mathbf{c}^*} |t_j| f_j(\mathbf{u})}{|\{\mathbf{x}^{(i)} : \mathbf{c}^*(\mathbf{x}^{(i)}) = 1\}|} \leq \mathcal{O}(t(\epsilon + \gamma) H(\mu N)).$$

Since every iteration covers at least  $\frac{\mu\gamma}{3t} N$  additional points,  $\hat{\mathbf{c}}$  has at most  $\frac{3}{\gamma} t$  terms. Thus the above bounds for small formulas ensure that  $\hat{\mathbf{c}}$  indeed has  $\Pr[\hat{\mathbf{c}}(\mathbf{x}) = 1] \geq (1 - \gamma)\mu$  and, by our first observation, the loss of  $\hat{\mathbf{c}}$  is indeed at most  $\mathcal{O}(t(\epsilon + \gamma) H(\mu N))$  as claimed. ■

### 5.3 Generalization Bound

Theorem 4.3 guarantees that Algorithm 4 finds a parameter vector approximating the optimal empirical loss, and Lemma 5.2 gives us the sample complexity needed to find a good condition for a fixed linear predictor. However, there is still a gap between the empirical loss and true loss. In this section, we will bound the generalization error of linear regression on each possible  $k$ -DNF, and then take a union bound to achieve the main theorem. In short, the process will blow-up the complexity by  $d^3$ , where  $d$  is the dimension of the feature space.

We will use the Rademacher generalization bound for linear predictors. For a set of data, Lemma 5.3 bounds the gap between the expected loss  $L_p(\cdot)$  and the empirical loss  $\hat{L}_p(\cdot)$ :

**Lemma 5.3 (Bartlett & Mendelson (2002), Kakade et al. (2009))** *For  $b > 0$ ,  $p \geq 1$ , random variables  $(\mathbf{Y}, Z)$  distributed over  $\{\mathbf{y} \in \mathbb{R}^d : \|\mathbf{y}\|_2 \leq b\} \times [b, b]$ , and any  $\delta \in (0, 1)$ , let  $L_p(\mathbf{w})$  denote  $\mathbb{E}[|\langle \mathbf{w}, \mathbf{Y} \rangle - Z|^p]$ , and for an i.i.d. sample of size  $N$  let  $\hat{L}_p(\mathbf{w})$  be the empirical loss  $\frac{1}{N} \sum_{j=1}^N |\langle \mathbf{w}, \mathbf{y}^{(j)} \rangle - z^{(j)}|^p$ . We then have that with probability  $1 - \delta$  for all  $\mathbf{w}$  with  $\|\mathbf{w}\|_2 \leq b$ ,*

$$|L_p(\mathbf{w}) - \hat{L}_p(\mathbf{w})| \leq \frac{2pb^{p+1}}{\sqrt{N}} + b^p \sqrt{\frac{2 \ln(4/\delta)}{N}}.$$

In our case, we only consider squared error; in other words,  $p = 2$  for us. And notice, in our setting, we are given a bound  $B$  on the magnitude of the entries, so  $b \leq \sqrt{dB}$ . Equivalently, we get

$$|L_p(\mathbf{w}) - \hat{L}_p(\mathbf{w})| \leq \frac{4B^3 d^{\frac{3}{2}}}{\sqrt{N}} + o(B^2 d).$$

Therefore, to bound the gap of the expected loss  $L_p(\cdot)$  and the empirical loss  $\hat{L}_p(\cdot)$ , it suffices for our sample complexity  $N$  to grow with  $B^6 d^3$ .

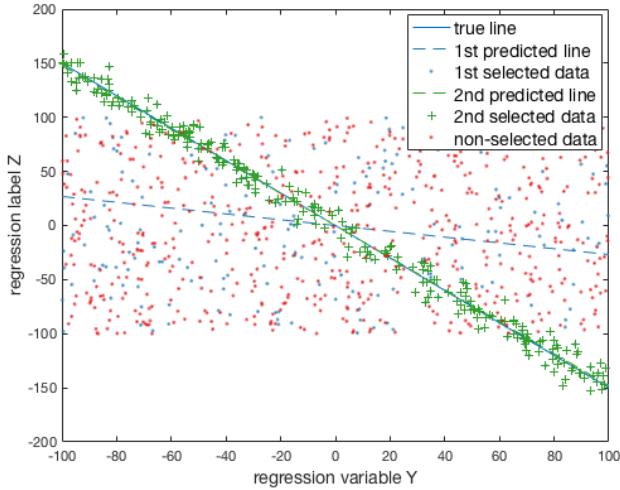
The above lemma bounds the gap of a specific (conditional) distribution and set of data. To obtain a bound on the loss obtained by our algorithm in general, we can simply take a union bound over the conditional distributions given by all  $t$ -term  $k$ -DNFs. Since  $x$  has  $n$  attributes, there are  $\binom{n}{k}$  terms, which is at most  $m = n^k$ . And, there are  $\binom{m}{t}$   $t$ -term  $k$ -DNFs, so in total we have  $\mathcal{O}(n^{kt})$  such  $k$ -DNFs, which means it will suffice to replace  $\delta$  with  $\frac{\delta}{n^{kt}}$  before taking the union bound. Actually, recalling that Lemma 5.2 only guarantees that our algorithm produces  $3t/\gamma$ -term  $k$ -DNFs as output, overall we thus achieve a  $\mathcal{O}(t \log(\mu N)(\gamma + \epsilon))$  approximation as claimed in the main theorem with  $N = \mathcal{O}(\frac{B^6 d^3 \sigma^2 L^2 t}{\mu \gamma^3} \log(\frac{m}{\delta/m^{t/\gamma}})) = \mathcal{O}(\frac{B^6 d^3 \sigma^2 L^2 t^2}{\mu \gamma^4} \log(m/\delta))$  examples.

## 6 Experiments

We now present several experiments. We present three synthetic data experiments with planted solutions to illustrate our algorithm's capabilities. For these synthetic data experiments, we use our algorithm to generate a list of candidate parameters and their corresponding DNF. If there is one pair that is close to our planted solution (or another output with even lower error), then we view the task as successful. We then present experiments showing that on real data sets, our algorithm can obtain loss that is consistently similar to or significantly smaller than that of the sparse  $\ell_2$  regression algorithm (Hainline et al., 2019). The code is written in Matlab with the optimization toolbox Yalmip Löfberg (2004). Our code can be found at <https://github.com/wumming/lud>.

### 6.1 Toy Examples

We first present a couple of toy examples that provide nice visual illustrations of our task.



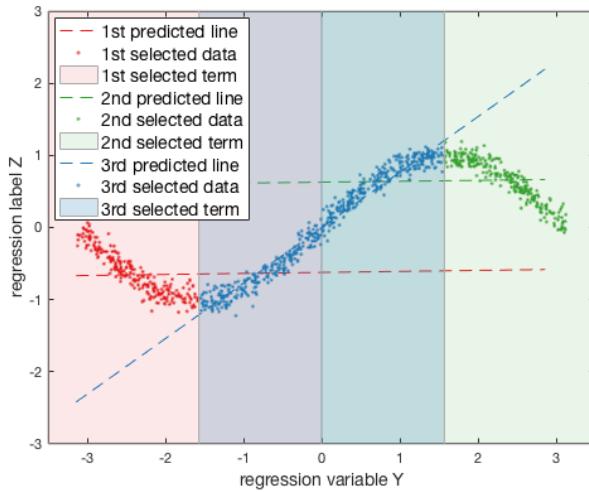
**Figure 4:** Line with Uniform

Conditional linear regression with  $d = 1$  in the  $y \times z$  plane: For the bad data,  $y, z \in [-100, 100]$ , while for good data,  $y \in [-100, 100]$  and  $z = -1.5y + \text{noise}$ , the  $\dim(\mathbf{x}) = 6$ . The algorithm gives two pairs of candidates. The blue one approximates the uniformly generated bad data and the green one catches the planted good data.

**Example 1: Line with Uniform** In this experiment, we first choose a 4-term 2-DNF at random. We uniformly generate Boolean attributes serving as  $\mathbf{x}$ , where  $\mu N$  of them satisfy the chosen DNF (good data) and the other  $(1 - \mu)N$  don't (bad data). Then the  $\mathbf{y}$  parts are all uniformly generated real attributes. We also generate a target linear rule  $\mathbf{w}^*$  with dimension equal to that of  $\mathbf{y}$ . For the good data, we set their labels  $z^{(i)} = \langle \mathbf{y}^{(i)}, \mathbf{w}^* \rangle + \text{noise}$ , where the noise is independently generated from zero-mean Gaussian distribution. For the bad data,  $z^{(i)}$  is independently generated from a uniform distribution similar to  $\mathbf{y}$ . We set the dimension of  $\mathbf{y}$  simply to be 1, as an illustration on the  $y \times z$  plane. There are  $N = 1000$  points in total and  $\mu = 0.25$ . For the bad data  $y^{(i)}, z^{(i)} \in [-100, 100]$ , while for good data  $y^{(i)} \in [-100, 100]$  and  $z^{(i)} = -1.5y^{(i)} + \text{noise}$ . The dimension of  $\mathbf{x}$  is set to 6.

Consequently there are  $m = 72$  terms in total. As shown in Figure 4, the algorithm finds two pairs of potential linear predictors and DNFs, of which the green one overlaps with the chosen planted rule, and the blue one matches the uniformly generated bad points.

**Example 2: Sine function** In this experiment, we didn't specify good and bad data, but uniformly generated  $y^{(i)} \in [-\pi, \pi]$  and  $z^{(i)} = \sin(y^{(i)}) + noise$ . We attach a constant coordinate to  $y$  to represent the intercept of the line. Define  $x_1, x_2$  and  $x_3$  to be “ $y \geq -\pi/2$ ”, “ $y \geq 0$ ” and “ $y \geq \pi/2$ ” respectively and set  $\mu = 0.5$  and  $S = 0.01$ . That means, the algorithm is asked to find a DNF segment that contains at least 50% of points that can fit a line. The algorithm outputs three pairs as shown on the Figure 5. The blue cluster  $x_1 \wedge \neg x_3$  has the smallest error, which is the points in interval  $[-\pi/2, \pi/2]$ .



**Figure 5:** Sine Function

Conditional linear regression with  $d = 1$  in the  $y \times z$  plane: for all data,  $y \in [-\pi, \pi]$ ,  $z = \sin(y) + noise$ ,  $dim(\mathbf{x}) = 3$ .  $x_1, x_2, x_3$  are defined as “ $y \leq -\pi/2, 0, \pi/2$ ” respectively. The algorithm gives three candidates: the red one is the region  $y < 0$ ; the green one  $y \leq 0$ ; the blue one  $-\pi/2 \leq y \leq \pi/2$ .

## 6.2 Example 3: Scale Up

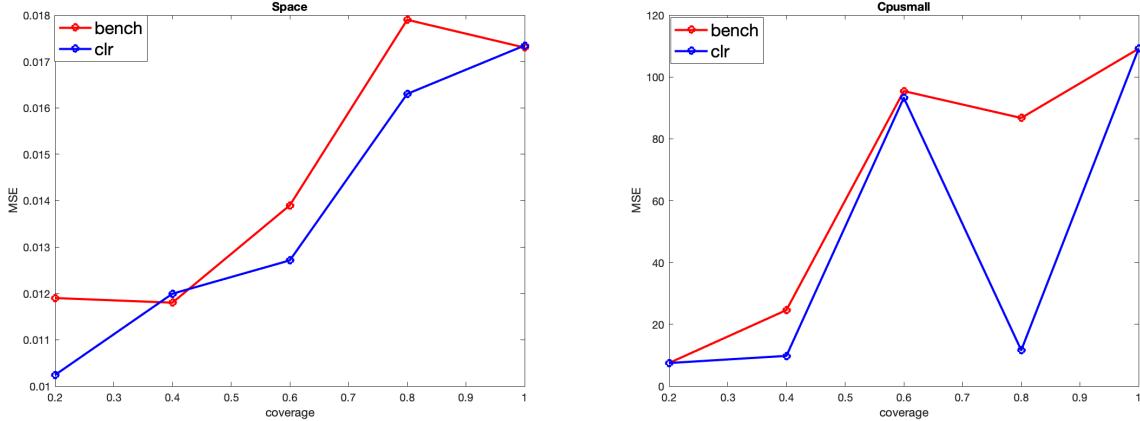
In this synthetic data experiment, we set  $dim(\mathbf{x}) = 7$ ,  $dim(\mathbf{y}) = 10$ , and generated  $N = 100000$  points in total with  $\mu = 0.5$ . We generate  $\mathbf{w}^*$  uniformly from  $[-10, 10]$  and randomly choose a 4-term 2-DNF to use as the condition describing the good data. For the bad data,  $\mathbf{y}^{(i)} \in [-1, 1]$ ,  $z^{(i)} \in [-10, 10]$ . For good data,  $\mathbf{y}^{(i)} \in [-1, 1]$  and  $z^{(i)} = \langle \mathbf{y}^{(i)}, \mathbf{w}^* \rangle + noise$  with variance 100. We set  $S = 0.1$ ,  $\gamma = 0.1$ ,  $r = 100$  and  $r_{final} = 1$ . In several trials, each time our algorithm finds several pairs of regression parameters and DNFs, with one of them containing our planted DNF. We note that there are also other pairs in the output with even lower error, meaning the algorithm also finds other subsets on which the error is smaller than for our planted solution. Note that since the previous algorithms (Juba, 2017; Hainline et al., 2019) scale exponentially with  $d$ , such an instance would be infeasible to solve using those algorithms.

### 6.3 Example 4: Real World Data

We test our algorithm on two of the larger benchmark data sets from the LIBSVM repository Chang & Lin (2011), Space and Cpusmall, used previously by Hainline et al. (2019), and compared the loss achieved for several target fractions. These standard data sets contains only real-valued attributes. Following Hainline et al's strategy, we generate Boolean attributes using indicators for membership in the empirical 50%-quantile of each real attributes. In other words, for each  $i$ th point, its  $j$ th Boolean attribute is defined as

$$x_j^{(i)} = \begin{cases} 1 & \text{if } y_j^{(i)} \geq \text{median}(y_j) \\ 0 & \text{if } y_j^{(i)} < \text{median}(y_j) \end{cases}$$

We randomly selected 1/3 of the data as the training data and the other 2/3 as the testing data. Similar to Hainline et al. (2019), we use the algorithm to find a (list of) DNF on the training data, compute the linear regression rule on the corresponding subset on the training data, and then test the loss of the regression rule on the corresponding subset of the testing data. Partially to cope with instability we observed in the SDP solver, we decreased the number of padded decompositions (in Algorithm 4) from  $112 \log(\frac{l(l+1)}{\delta})$  to 12, and instead repeated the algorithm 50 times to produce a list of candidates. We compute the linear fit and training loss for each candidate DNF, and then use the DNF with lowest training loss as our final output to test on the testing data. We repeated the experiments using different values of  $\mu$  (0.2, 0.4, 0.6, 0.8, 1.0) and set the parameters  $S = \mu \cdot 10$ ,  $\gamma = 0.1$ , and  $r_{final} = S \cdot \dim(y)$ , where  $\dim(y)$  is the dimension of the real features.



**Figure 6:** Linear regression on LIBSVM datasets  
**Left: Space data:**  $N = 1025$ ,  $\dim(x) = \dim(y) = 6$ . **Right: Cpusmall data:**  $N = 2703$ ,  $\dim(x) = \dim(y) = 12$ .

As shown in Figure 6, our performance is better than benchmark algorithm for most of the cases, and at least comparable in all cases. For  $\mu = 1$ , conditional linear regression is the same as a simple linear regression, so unsurprisingly, our algorithm and Hainline et al's algorithm have the same result. We note that our running time is much better than Hainline et al's algorithm, even on these small data sets. Their algorithm required a few days of cloud computing time, while our algorithm only required a few hours.

## 6.4 Discussion of the Experiments

These experiments illustrate the variety of tasks our algorithm can solve. They also demonstrate that our algorithm can work in practice: The algorithm is feasible to run, in spite of the fact that it must solve numerous SDP optimization problems. It does find the desired, planted answers on the synthetic data where this is known, and it obtains comparable or superior error to the baseline algorithm of Hainline et al. (2019). We obtain a computational advantage over Charikar et al. (2017) because we group the data points into terms and precompute the loss functions for these terms. So, instead of solving a huge SDP problem with variables for each example as in Charikar et al’s approach, our SDPs are potentially of reasonable size and can run in a fraction of a second. The radius  $r$  of the candidate parameters shrinks exponentially over the outer loop of the list regression algorithm, so even if we set  $r_{final}$  to be 0.1 or so, it terminates within 10 iterations. These improvements yield an algorithm that is feasible to run in practice, so long as the number of regression features and candidate terms is only moderately large.

One difficulty with the use of our algorithm is that we need to specify the setting of several parameters. One has to guess an estimate of  $\mu$  and  $S_0$  to run the algorithm, so in practice one may need to try a sequence of candidate settings for these parameters. The smaller  $S$  and  $r_{final}$  are set, the more accurate an estimate the algorithm can provide, at the cost of more computation time. In practice, we cannot set  $r_{final}$  to be too small, as otherwise the algorithm tends to think there does not exist a candidate with size  $\mu$ . We believe that this may be a consequence of noise in the data.

It is interesting that, on the Cpusmall dataset, the loss obtained at  $\mu = 0.8$  is much smaller than the loss we obtained at  $\mu = 0.6$ . Usually we would expect the loss will decrease with  $\mu$ : since our problem statement allows us to output a solution on which the condition comprises 80% of the data when we are only seeking 60%, it is strictly easier to fit a smaller subset. We note that this difference is not caused by the double-counting of points, since the double-counted loss of  $\mu = 0.8$  is still smaller than the loss of  $\mu = 0.6$ . Thus, it seems that the problem is that the algorithm does not obtain such accurate estimates of the regression parameters on  $\mu = 0.6$ . We note that the value of  $\mu$  enters Algorithm 1 in several places, and the guarantees we can provide on its output feature a  $1/\sqrt{\mu}$  factor blow-up in the error; notice, our Theorem 2.3 requires a stronger condition for smaller  $\mu$ ). Indeed, in the soft regression/outlier detection procedure, we might intuitively expect that as the proportion of “signal” decreases, we might get less accurate estimates. Thus, in practice, we would advise running the algorithm several times using  $\mu = (1 - \Delta), (1 - \Delta)^2, \dots, (1 - \Delta)^i$  for some number of iterations  $i$  ( $\Delta \in (0, 1)$ ) down to a desired minimum, and taking the condition/parameters with the smallest loss.

One may worry that for small  $\mu$ , there is a chance of overfitting. However, this is not a problem as long as we have a reasonably large data set and enough runs. We observe that for  $\mu = 0.2$ , 10 out of the 50 runs do have much larger testing error than training error, which implies overfitting. But when we pick the candidate of the minimum training error, overfitting doesn’t impact our output. To prevent overfitting, we suggest to use larger data sets and run more trials so that even small subsets have lower variance.

We also caution that our algorithm requires enough data to achieve stability. We also examined some of the smaller LIBSVM data sets (bodyfat and Boston housing) that have only a few hundred examples. For these small data sets, we found that the SDP solver could not solve instances in which the radius of the parameter space decreased below 1. Thus, the final iteration of parameters we obtained provided poor estimates that were not competitive with the algorithm of Hainline et al. (2019). In conclusion, our algorithm is better in terms of both running time and accuracy than

the previous algorithm for the larger benchmark data sets ( $N \geq 1000$ ). Since it collapses the data points into loss matrices for terms, and the corresponding pre- and post-processing can be done in roughly linear time, it can scale up to very large data sets. But, it is not appropriate to use for very small data sets.

## 7 Directions for future work

Our work has introduced the first algorithm for conditional linear regression that simultaneously obtains conditions of near-optimal probability, obtains an  $\mathcal{O}(t \log \mu N)$ -approximation to the optimal loss, and does not feature a running time that depends exponentially on the number of factors used in the linear predictor. Theoretically, the one “additional cost” imposed by our approach was a requirement that  $S_{\varepsilon 0}$ , the spectral difference between the covariances of the individual terms and the overall desired subset of the data we are interested in, is sufficiently small compared to the convexity of our loss functions. Although we noted that we can always impose sufficient convexity by adding a regularization term, in some cases this may yield insufficiently accurate estimates of the regression parameters. Thus, the main theoretical question raised here is whether or not this requirement can be relaxed or removed entirely. We stress that the previous approaches did not require such an assumption, which suggests that it may be possible to remove it.

Another family of questions concerns the amount of blow-up we incur over the loss. While our bound seems quite good when the number of terms in the desired DNF  $t$  is small, one can always ask whether or not it is optimal. We do not have techniques for answering this question at the moment. At the same time, we recall that Zhang et al. (2017) obtain a better approximation factor for the condition search (learning abduction) problem for large formulas: by more carefully controlling the amount of double-counting, Zhang et al. obtain a better blow-up when  $t \gg n^{k/2}$  for  $k$ -DNFs on  $n$  Boolean attributes. Thus, another natural question is if we can likewise quadratically reduce the blow-up of the loss, as achieved by Zhang et al. for the condition search problem.

One can also seek to further improve the efficiency of our algorithm in several respects. For one, we have made no attempt to optimize the amount of data required for our guarantees. (Or, for that matter, required by the algorithm.) For two, we believe that there is substantial scope for improving the running time of our algorithm. Diakonikolas et al. (2018) recently made several improvements to the algorithm of Charikar et al. (2017), replacing the semidefinite programs with eigenvector computations and using a simpler clustering strategy than padded decompositions. It may be possible to analogously improve the running time of our algorithm.

Finally, we note that Charikar et al. (2017) had intended their algorithm for use in solving a variety of tasks, and thus gave a generic analysis in terms of abstract loss functions. This helped enable our work, although our results depend on key features of the squared-error loss for regression—in particular, we needed that we could show that our loss decreased linearly with the radius. Then, since we could give a bound on the radius that scaled linearly with the previous iteration’s radius, we could iteratively improve our estimates by simply running another iteration of the algorithm. There are conceptually similar tasks, such as conditional *linear classification*, where we seek to make Boolean (as opposed to real-valued) predictions, that we might hope to be able to solve using an algorithm of this form, but our approach would require that linear scaling condition holds for the classification loss function, which is not clear. It would be interesting to see whether and how our approach could be adapted to solve such problems.

## Acknowledgements

We thank Jacob Steinhardt for sharing his code and many helpful discussions. We also thank Ilias Diakonikolas for a helpful discussion. We finally thank our anonymous reviewers for their comments on related works.

## References

- Bartlett, P. L. and Mendelson, S. Rademacher and gaussian complexities: Risk bounds and structural results. *Journal of Machine Learning Research*, 3(Nov):463–482, 2002.
- Chang, C.-C. and Lin, C.-J. Libsvm: a library for support vector machines. *ACM transactions on intelligent systems and technology (TIST)*, 2(3):27, 2011.
- Charikar, M., Steinhardt, J., and Valiant, G. Learning from untrusted data. In *Proceedings of the 49th Annual ACM SIGACT Symposium on Theory of Computing, STOC 2017, Montreal, QC, Canada, June 19-23, 2017*, pp. 47–60, 2017. doi: 10.1145/3055399.3055491. Full version arXiv:1611.02315v2 [cs.LG].
- Cortes, C., DeSalvo, G., and Mohri, M. Learning with rejection. In *ALT 2016*, volume 9925 of *LNAI*, pp. 67–82. 2016.
- Diakonikolas, I., Kane, D. M., and Stewart, A. List-decodable robust mean estimation and learning mixtures of spherical Gaussians. In *Proceedings of the 50th Annual ACM SIGACT Symposium on Theory of Computing*, pp. 1047–1060, 2018.
- El-Yaniv, R. and Wiener, Y. Pointwise tracking the optimal regression function. In *Advances in Neural Information Processing Systems*, pp. 2042–2050, 2012.
- Fakcharoenphol, J., Rao, S., and Talwar, K. A tight bound on approximating arbitrary metrics by tree metrics. In Larmore, L. L. and Goemans, M. X. (eds.), *Proceedings of the 35th Annual ACM Symposium on Theory of Computing, June 9-11, 2003, San Diego, CA, USA*, pp. 448–455. ACM, 2003. ISBN 1-58113-674-9. doi: 10.1145/780542.780608. URL <http://doi.acm.org/10.1145/780542.780608>.
- Fischler, M. A. and Bolles, R. C. Random sample consensus: A paradigm for model fitting with applications to image analysis and automated cartography. *Communications of the ACM*, 24(6):381–395, 1981.
- Fröhlich, B., Rodner, E., Kemmler, M., and Denzler, J. Large-scale gaussian process classification using random decision forests. *Pattern Recognition and Image Analysis*, 22, 2012. URL <https://link.springer.com/article/10.1134%2FS1054661812010166>.
- Hainline, J., Juba, B., Le, H. S., and Woodruff, D. P. Conditional sparse  $\ell_p$  regression with optimal probability. In *Proc. 22nd AISTATS*, volume 89 of *PMLR*, pp. 369–382, 2019.
- Haussler, D. Quantifying inductive bias: AI learning algorithms and Valiant’s learning framework. *Artificial Intelligence*, 36:177–221, 1988.

- Huber, P. J. *Robust Statistics*. John Wiley & Sons, New York, NY, 1981.
- Jiang, J. *Linear and Generalized Linear Mixed Models and Their Applications*. Springer, Berlin, 2007.
- Juba, B. Conditional sparse linear regression. In *8th Innovations in Theoretical Computer Science Conference, ITCS 2017, January 9-11, 2017, Berkeley, CA, USA*, pp. 45:1–45:14, 2017. doi: 10.4230/LIPIcs.ITCS.2017.45. URL <https://doi.org/10.4230/LIPIcs.ITCS.2017.45>.
- Juba, B., Li, Z., and Miller, E. Learning abduction under partial observability. In *Proceedings of the Thirty-Second AAAI Conference on Artificial Intelligence, New Orleans, Louisiana, USA, February 2-7, 2018*, 2018. URL <https://www.aaai.org/ocs/index.php/AAAI/AAAI18/paper/view/17323>.
- Kakade, S. M., Sridharan, K., and Tewari, A. On the complexity of linear prediction: Risk bounds, margin bounds, and regularization. In *Advances in neural information processing systems*, pp. 793–800, 2009.
- Löfberg, J. Yalmip : A toolbox for modeling and optimization in matlab. In *In Proceedings of the CACSD Conference*, Taipei, Taiwan, 2004.
- McCulloch, C. E. and Searle, S. R. *Generalized, Linear, and Mixed Models*. John Wiley & Sons, New York, NY, 2001.
- Park, Y. W., Jiang, Y., Klabjan, D., and Williams, L. Algorithms for generalized cluster-wise linear regression. *INFORMS Journal on Computing*, 29(2):301–317, 2017.
- Quinlan, J. R. Learning with continuous classes. In *5th Australian joint conference on artificial intelligence*, volume 92, pp. 343–348. Singapore, 1992.
- Rousseeuw, P. J. and Leroy, A. M. *Robust Regression and Outlier Detection*. John Wiley & Sons, New York, NY, 1987.
- Slavík, P. Improved performance of the greedy cover algorithm for partial cover. *Information Processing Letters*, 64(5):251–254, 1997.
- Zhang, M., Mathew, T., and Juba, B. An improved algorithm for learning to perform exception-tolerant abduction. In *Proc. 31st AAAI*, pp. 1257–1265, 2017.