# Sparsifying Generalized Linear Models\*

# Arun Jambulapati

Simons Institute for the Theory of Computing Berkeley, USA jmblpati@gmail.com

# Yang P. Liu

Institute for Advanced Study Princeton, USA yangpatil@gmail.com

#### **ABSTRACT**

We consider the sparsification of sums  $F:\mathbb{R}^n\to\mathbb{R}^+$  where  $F(x)=f_1(\langle a_1,x\rangle)+\cdots+f_m(\langle a_m,x\rangle)$  for vectors  $a_1,\ldots,a_m\in\mathbb{R}^n$  and functions  $f_1,\ldots,f_m:\mathbb{R}\to\mathbb{R}^+$ . We show that  $(1+\varepsilon)$ -approximate sparsifiers of F with support size  $\frac{n}{\varepsilon^2}(\log\frac{n}{\varepsilon})^{O(1)}$  exist whenever the functions  $f_1,\ldots,f_m$  are symmetric, monotone, and satisfy natural growth bounds. Additionally, we give efficient algorithms to compute such a sparsifier assuming each  $f_i$  can be evaluated efficiently.

Our results generalize the classical case of  $\ell_p$  sparsification, where  $f_i(z) = |z|^p$ , for  $p \in (0,2]$ , and give the first near-linear size sparsifiers in the well-studied setting of the Huber loss function and its generalizations, e.g.,  $f_i(z) = \min\{|z|^p, |z|^2\}$  for  $0 . Our sparsification algorithm can be applied to give near-optimal reductions for optimizing a variety of generalized linear models including <math>\ell_p$  regression for  $p \in (1,2]$  to high accuracy, via solving  $(\log n)^{O(1)}$  sparse regression instances with  $m \le n(\log n)^{O(1)}$ , plus runtime proportional to the number of nonzero entries in the vectors  $a_1, \ldots, a_m$ .

### **CCS CONCEPTS**

• Theory of computation → Sketching and sampling.

# **KEYWORDS**

Generalized linear models, sparsification, chaining, Lewis weights

#### **ACM Reference Format:**

Arun Jambulapati, James R. Lee, Yang P. Liu, and Aaron Sidford. 2024. Sparsifying Generalized Linear Models. In *Proceedings of the 56th Annual ACM Symposium on Theory of Computing (STOC '24), June 24–28, 2024, Vancouver, BC, Canada.* ACM, New York, NY, USA, 11 pages. https://doi.org/10.1145/3618260.3649684

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

STOC '24, June 24-28, 2024, Vancouver, BC, Canada

© 2024 Copyright held by the owner/author(s). Publication rights licensed to ACM. ACM ISBN 979-8-4007-0383-6/24/06

https://doi.org/10.1145/3618260.3649684

#### James R. Lee

University of Washington Seattle, USA jrl@cs.washington.edu

#### Aaron Sidford

Stanford University Stanford, USA sidford@stanford.edu

### 1 INTRODUCTION

Empirical risk minimization (ERM) is a widely studied problem in learning theory and statistics (see, e.g., [18], for relevant references to the expansive literature on this topic). A prominent special case is the problem of optimizing a *generalized linear model (GLM)*, i.e.,

$$\min_{x \in \mathbb{R}^n} F(x) \quad \text{for} \quad F(x) := \sum_{i=1}^m f_i(\langle a_i, x \rangle - b_i), \quad (1.1)$$

where the total loss  $F: \mathbb{R}^n \to \mathbb{R}$ , is defined by vectors  $a_1, \ldots, a_m \in \mathbb{R}^n$ ,  $b \in \mathbb{R}^m$ , and loss functions  $f_1, \ldots, f_m : \mathbb{R} \to \mathbb{R}$ . Different choices of the loss functions  $\{f_i\}$  capture important problems, including linear regression, logistic regression, and  $\ell_p$  regression [2, 8].

Recently, efficient algorithms for solving (1.1) to high-accuracy have been developed in many settings [5, 7, 13] such as linear programming and  $\ell_1$ -regression, where  $f_i(x) = |x|$ . For example, when m is on the order of n, it is known how to solve linear programs and some GLMs in roughly (up to logarithmic factors) the time it currently takes to multiply two general  $n \times n$  matrices [2, 11, 16, 18] which is, up to logarithmic factors, the best-known, running time for solving a single linear system in a dense  $n \times n$  matrix.

When  $m \gg n$ , a natural approach for fast algorithms is to apply sparsification techniques to reduce the value of m, while maintaining a good multiplicative approximation of the objective value. More precisely, say that the objective F admits an s-sparse  $\varepsilon$ -approximation if there are non-negative weights  $w_1, \ldots, w_m \in \mathbb{R}_+^m$ , at most s of which are non-zero, and such that

$$|F(x) - \tilde{F}(x)| \le \varepsilon F(x)$$
 for all  $x \in \mathbb{R}^n$ , where 
$$\tilde{F}(x) := \sum_{i=1}^m w_i f_i(\langle a_i, x \rangle - b_i).$$

When  $f_i(z)=|z|^p$  are  $\ell_p$  losses, near-optimal sparsification results are known: If p>0, then F admits an s-sparse  $\varepsilon$ -approximation for  $s \leq \tilde{O}(n^{\max\{1,p/2\}}\varepsilon^{-2});^1$  this sparsity bound is known to be optimal up to polylogarithmic factors [4, 22, 25, 26]. In particular, for  $p\in(0,2]$ , the size is  $\tilde{O}(n\varepsilon^{-2})$ , near-linear in the underlying dimension n. The p=2 case has been especially influential in the development of several fast algorithms for linear programming and graph optimization over the last two decades [6, 23, 24].

 $<sup>^{\</sup>star}\mathrm{A}$  full version of the paper can be found at https://arxiv.org/pdf/2311.18145.pdf.

 $<sup>^{\</sup>overline{1}}$  Throughout, we use  $\tilde{O}(f)$  to suppress polylogarithmic quantities in  $\emph{m},\emph{n},\varepsilon^{-1},$  and f

However, as far as the authors know,  $\ell_p$  losses are the only class of natural loss functions for which linear-size sparsification results are known for GLMs. For instance, for the widely-studied class of Huber loss functions (see (1.4)) and related variants, e.g.,  $f_i(z) = \min\{|z|, |z|^2\}$ , the best known sparsity bound was  $\tilde{O}(n^{4-2\sqrt{2}}\varepsilon^{-2})$  [20]. Improving this bound to near-linear (in n) is an established important open problem that has potential applications to regression for Huber and  $\ell_p$  losses [1, 3, 13, 20, 28].

The main result of this paper is near-optimal sparsification for a large family of loss functions  $\{f_i\}$  that include the Huber losses,  $\ell_p$  losses, and generalizations. Informally, we show that if the loss functions  $\{f_i\}$  are nonnegative, symmetric, and grow at most quadratically, then there exists an s-sparse  $\varepsilon$ -approximation of F with  $s \leq \tilde{O}(n\varepsilon^{-2})$ . Moreover, the sparse approximation can be found very efficiently, in time proportional to the time used for  $\tilde{O}(1)$  instances of  $\ell_2$ -sparsification (Theorem 1.1). A particularly nice application of our result is an algorithm that solves  $\ell_p$ -regression to high accuracy for  $1 by reducing to <math>\tilde{O}_p(1)$  instances of  $\ell_p$ -regression with  $m = \tilde{O}(n)$  (Theorem 1.2). Our framework can also be applied to minimizing sums of  $\gamma_p$  functions for  $p \in (1,2]$  (see (1.4)) to high accuracy, and to approximate Huber regression.

The main technical hurdle in obtaining these results is that the loss functions are not necessarily homogeneous and they can exhibit different behaviors at different scales. Note that this hurdle arises already for losses like  $f_i(z) = \min\{|z|, |z|^2\}$ , even though the loss function only has two different scaling regimes. To overcome this hurdle we develop a multiscale notion of "importance scores" for appropriately down-sampling F into a sparse representation.

# 1.1 Hypotheses and Results for Sparsification

Consider a generalized linear model as in (1.1), with loss functions  $f_1, \ldots, f_m : \mathbb{R} \to \mathbb{R}_+$  and vectors  $a_1, \ldots, a_m \in \mathbb{R}^n$ . For simplicity, we assume that b = 0 in (1.1). This is without loss of generality, as  $\langle a_i, x \rangle - b_i = \langle (a_i, b_i), (x, -1) \rangle$ , and  $(a_i, b_i), (x, -1) \in \mathbb{R}^{n+1}$ , so we can re-encode the problem in n + 1 dimensions with b = 0.

We will often think of the case  $f_i(z) = h_i(z)^2$  for some  $h_i : \mathbb{R} \to \mathbb{R}_+$ , as the assumptions we need are stated more naturally in terms of  $\sqrt{f_i}$ . To that end, consider a function  $h : \mathbb{R}^k \to \mathbb{R}_+$  and the following two properties, where  $L \ge 1$  and  $c, \theta > 0$  are some positive constants.<sup>2</sup>

- (a) (*L*-auto-Lipschitz)  $|h(z) h(z')| \le L h(z z')$  for all  $z, z' \in \mathbb{R}^k$ . (b) (Lower  $\theta$ -homogeneous)  $h(\lambda z) \ge c \lambda^{\theta} h(z)$  for all  $z \in \mathbb{R}^k$  and
- Note that if  $h : \mathbb{R} \to \mathbb{R}$  is concave and symmetric, then it is 1-auto-Lipschitz.

We can now state our main theorem, whose proof appears in the full version.

Theorem 1.1. Consider  $f_1, \ldots, f_m : \mathbb{R} \to \mathbb{R}_+$ , and suppose there are numbers  $L \geq 1$ ,  $c, \theta > 0$  such that each  $\sqrt{f_i}$  is L-auto-Lipschitz and lower  $\theta$ -homogeneous (with constant c). Then for any  $a_1, \ldots, a_m \in \mathbb{R}^n$ , and numbers  $0 < \varepsilon < \frac{1}{2}$  and  $s_{\max} > s_{\min} \geq 0$ , there are

nonnegative weights  $w_1, \ldots, w_m \ge 0$  such that

$$\left| F(x) - \sum_{i=1}^{m} w_i f_i(\langle a_i, x \rangle) \right| \le \varepsilon F(x),$$

$$\forall x \in \mathbb{R}^n \text{ s.t. } s_{\min} \le F(x) \le s_{\max},$$

where  $F(x) := f_1(\langle a_1, x \rangle) + \cdots + f_m(\langle a_m, x \rangle)$ , and

$$|i \in \{1, \dots, m\} : w_i > 0| \lesssim_{L,c,\theta} \frac{n}{\varepsilon^2} \log \left(\frac{n}{\varepsilon} \frac{s_{\max}}{s_{\min}}\right) (\log S)^3, \quad (1.2)$$

where 
$$S := \frac{n}{\varepsilon} \log \left( \frac{2s_{\text{max}}}{s_{\text{min}}} \right)$$
. (1.3)

Moreover, with high probability, the weights  $\{w_i\}$  can be computed in time

$$\tilde{O}_{L,c,\theta}((\operatorname{nnz}(a_1,\ldots,a_m)+n^\omega+m\mathcal{T}_{\operatorname{eval}})\log(ms_{\max}/s_{\min}))$$
.

Here,  $\mathcal{T}_{\text{eval}}$  is the maximum time needed to evaluate each  $f_l$ ,  $\operatorname{nnz}(a_1,\ldots,a_m)$  is the total number of non-zero entries in the vectors  $a_1,\ldots,a_m$ , and  $\omega$  is the matrix multiplication exponent. "High probability" means that the failure probability can be made less than  $n^{-\ell}$  for any  $\ell>1$  by increasing the running time by an  $O(\ell)$  factor.

We use the notation  $O_{L,c,\theta}$  and  $\tilde{O}_{L,c,\theta}(\cdot)$  to indicate an implicit dependence on the parameters  $L,c,\theta$ , and  $A \lesssim_{L,c,\theta} B$  is shorthand for  $A \leqslant O_{L,c,\theta}(B)$ . The constant hidden by the  $O_{L,c,\theta}(\cdot)$  notation is about  $(L/c)^{O(\theta^{-2})}$ , though we made no significant effort to optimize this dependence.

It is not difficult to see that for  $0 , the function <math>f_i(z) = |z|^p$  satisfies the required hypotheses of Theorem 1.1. In the full version, we show that  $\gamma_p$  functions, defined as

$$\gamma_{p}(z) := \begin{cases} \frac{p}{2}z^{2} & \text{for } |z| \leq 1\\ |z|^{p} - (1 - \frac{p}{2}) & \text{for } |z| \geq 1, \end{cases}$$
 (1.4)

for  $p \in (0, 2]$ , also satisfy the conditions. The special case of  $\gamma_1$  is known as the Huber loss. We note that this function can be generalized to other thresholds.

The  $\gamma_p$  functions were introduced in [8] and have since been used in several works on high-accuracy  $\ell_p$  regression [1, 2, 13]. Due to these connections, the works [13, 20] studied sparsification with  $\gamma_p$  losses, providing sparsity bounds of  $\tilde{O}(n^3)$  and  $\tilde{O}(n^{4-2\sqrt{2}}) \approx \tilde{O}(n^{1.172})$ , respectively. More precisely, [20] establish a bound of  $\tilde{O}(n^{1+\delta(p)})$  for  $p \in [1,2]$  with  $\delta(1) = 3 - 2\sqrt{2}$ , and  $\delta(p) \to 0$  as  $p \to 2$ .

# 1.2 Fast $\ell_p$ Regression

Combining our sparsification theorem with iterative refinement [2] yields near-optimal reductions for solving  $\ell_p$  regression to high accuracy. More specifically, we show that  $\ell_p$  regression for matrices  $A \in \mathbb{R}^{m \times n}$  can be reduced to a sequence of  $\tilde{O}_p(1)$  instances with  $\widetilde{A} \in \mathbb{R}^{\tilde{O}(n) \times n}$ . It is known how to solve such instances in time  $n^{\omega_0}$  for  $\omega_0 := 2 + \max\left\{\frac{1}{6}, \omega - 2, \frac{1-\alpha}{2}\right\}$  [18], where  $\alpha$  is the dual matrix multiplication exponent. Alternatively, they can each be solved in roughly  $n^{1/3}$  iterations and time  $n^{\max\{\omega, 2+1/3\}}$  [2], where an "iteration" refers to an operation that is dominated by the cost of solving a particular  $n \times n$  linear system.

<sup>&</sup>lt;sup>2</sup>The setting k = 1 suffices for the present work, though we state them for general k > 1

THEOREM 1.2 (FAST  $\ell_p$  REGRESSION). There is an algorithm that given any  $A \in \mathbb{R}^{m \times n}$ ,  $b \in \mathbb{R}^m$ , and  $p \in (1,2]$  computes an x satisfying

$$||Ax - b||_p^p \le (1 + \varepsilon) \min_{x \in \mathbb{R}^n} ||Ax - b||_p^p$$

in either  $\tilde{O}_p(n^{\frac{2-p}{p+2}})$  iterations and  $\tilde{O}_p(\operatorname{nnz}(A) + n^{\max\{\omega,2+1/3\}})$  time, or  $\tilde{O}_p(\sqrt{n})$  iterations and  $\tilde{O}_p(\operatorname{nnz}(A) + n^{\omega_0})$  time, with high probability.

It is standard to turn a high accuracy algorithm for an optimization problem into one that solves a corresponding dual problem. We present such an argument for  $\ell_p$ -regression in the full version.

Theorem 1.3 (Dual of  $\ell_p$  regression). There is an algorithm that given  $A \in \mathbb{R}^{m \times n}$ ,  $c \in \mathbb{R}^m$ , and  $q \in [2, \infty)$  computes  $a y \in \mathbb{R}^m$  satisfying  $A^\top y = c$  and

$$\|y\|_q^q \leqslant (1+\varepsilon) \min_{A^\top y = c} \|y\|_q^q$$

in either  $\tilde{O}_q(n^{\frac{q-2}{3q-2}})$  iterations and  $\tilde{O}_q(\text{nnz}(A) + n^{\max\{\omega,2+1/3\}})$  time, or  $\tilde{O}_q(\sqrt{n})$  iterations and  $\tilde{O}_q(\text{nnz}(A) + n^{\omega_0})$  time, with high probability.

Prior work [15] shows that  $\ell_p$ -regression can be solved in  $\tilde{O}_p(n^{1/3})$  iterations of solving a linear system for  $p \in [2, \infty)$ . Combining that result with Theorem 1.2 shows that  $\ell_p$ -regression can be solved using  $\tilde{O}_p(n^{1/3})$  linear systems for all p > 1.

# 1.3 Discussion of the Hypotheses

Let us now discuss various hypotheses and the extent to which they are necessary for sparsifiers of nearly-linear size to exist. In addition to the properties a and b, let us consider three others that we will use frequently. In what follows, C, u > 0 are positive constants and  $h : \mathbb{R}^n \to \mathbb{R}_+$ .

- (c) (*C*-symmetric)  $h(z) \leq Ch(-z)$  for all  $z \in \mathbb{R}^n$ .
- (d) (*C*-monotone)  $h(z) \le Ch(\lambda z)$  for  $\lambda \ge 1$ .
- (e) (Upper *u*-homogeneous)  $h(\lambda z) \leq C\lambda^u h(z)$  for all  $z \in \mathbb{R}^n$  and  $\lambda \geq 1$ .

First, note that a and b imply c-e.

LEMMA 1.4. The following implications hold:

- (1) h is L-auto-Lipschitz  $\implies h$  is L-symmetric.
- (2) h is lower  $\theta$ -homogeneous with constant  $c \implies h$  is 1/c-monotone
- (3) h is L-auto-Lipschitz and C-monotone ⇒ h is upper 1-homogeneous with constant 2CL.

PROOF. Since h(0)=0, applying the definition of L-auto-Lipschitz with z=0 gives  $h(-z)\leqslant Lh(z)$  for any  $z\in\mathbb{R}^n$ . The second implication is immediate. For the third, note that for a positive integer k, we have  $h(kz)\leqslant \sum_{j=0}^{k-1}|h((j+1)z)-h(jz)|\leqslant kLh(z)$ . Using the L-auto-Lipschitz property again gives

$$h(\lambda z) \leq h(\lceil \lambda \rceil z) + Lh((\lceil \lambda \rceil - \lambda)z)$$
  
$$\leq \lceil \lambda \rceil L \cdot h(z) + LCh(z) \leq 2CL\lambda h(z),$$

where the penultimate inequality uses *C*-monotonicity.

**Symmetry c.** To illustrate the need for approximate symmetry, let us consider gluing together two functions that are otherwise "nice" in our framework:

$$f(z) := \begin{cases} |z|^2, & z \geqslant 0 \\ |z|, & z < 0. \end{cases}$$

Suppose that  $f_1 = \cdots = f_m = f$ . Consider unit vectors  $\hat{a}_1, \ldots, \hat{a}_m \in \mathbb{R}^n$  such that  $\delta_{ij} \coloneqq |\langle \hat{a}_i, \hat{a}_j \rangle| < \frac{1}{2}$  for  $i \neq j$ . A basic volume computation shows that one can choose  $m \geqslant 2^{\Omega(n)}$ . Denote  $a_i \coloneqq (\hat{a}_i, 1) \in \mathbb{R}^{n+1}$  for  $i = 1, \ldots, m$ .

Then for  $\lambda > 0$  and  $x := \lambda(\hat{a}_i, -\frac{1}{2})$ , we have

$$\begin{split} f_{j}(\langle a_{j}, x \rangle) &= f(\langle a_{j}, \lambda(\hat{a}_{i}, -\frac{1}{2}) \rangle) \\ &= \begin{cases} f(\lambda/2) \times \lambda^{2} & i = j, \\ f(\lambda(\delta_{ij} - \frac{1}{2})) \lesssim \lambda & \text{otherwise.} \end{cases} \end{split}$$

Thus in any approximate sparsifier  $\tilde{F} = w_1 f_1 + \cdots + w_m f_m$ , it must be that either  $w_i > 0$ , or  $\sum_{j \neq i} w_j \gtrsim \lambda$ . Sending  $\lambda \to \infty$  shows that the latter is impossible.

**Lower growth and monotonicity b, d.** We consider these properties together since monotonicity is a weaker property than lower homogeneity. A natural function that does not satisfy lower homogeneity is the *Tukey loss* which, for the sake of the present discussion, one can take as  $f_i(z) := \min\{1, |z|^2\}$ , which is a natural analog of  $\gamma_p$  (recall (1.4)) for p = 0.

For sparsifying GLMs with the Tukey loss, previous works have made additional assumptions. For example, that one only ensures sparsification when  $\|a_i\|_2 \le n^{O(1)}$ , and for inputs  $x \in \mathbb{R}^n$  satisfying  $\|x\|_2 \le n^{O(1)}$ ; see [9, Assumption 2] and the discussion afterwards, and [20, §8.3]. In the full version, we show how to achieve a  $\tilde{O}(n^{1+o(1)}\varepsilon^{-2})$ -sparse  $\varepsilon$ -approximations under these assumptions. At a high level, the simple idea is to consider the proxy loss functions  $\hat{f}_i(z) := \min\{|z|^p, |z|^2\}$  with p sufficiently small.

**Upper quadratic growth e.** Note that, by Lemma 1.4, if  $\sqrt{f_i}$  satisfies a, then  $f_i$  is upper 2-homogeneous. For near-linear size sparsifiers, 2-homogeneity is a natural condition, since sparsifying with loss functions  $f_i(x) = |x|^p$  and p > 2 requires the sparsifier to have at least  $\Omega(n^{p/2})$  terms [4].

**The auto-Lipschitz property a.** As Lemma 1.4 shows, this property gives us approximate symmetry c and upper 1-homogeneity e. Crucially, this property also allows us to exploit the geometry of the vectors  $a_1, \ldots, a_m \in \mathbb{R}^n$ . Note that a implies

$$(f_i(z) - f_i(z'))^2 =$$

$$(f_i(z)^{1/2} - f_i(z')^{1/2})^2 (f_i(z)^{1/2} + f_i(z')^{1/2})^2$$

$$\leq 2L^2 f_i(z - z') (f_i(z) + f_i(z')).$$

In particular, we have

$$\left( f_i(\langle a_i, x \rangle) - f_i(\langle a_i, y \rangle) \right)^2$$

$$\leq 2L^2 \underbrace{f_i(\langle a_i, x - y \rangle)}_{} (f_i(\langle a_i, x \rangle + f_i(\langle a_i, y \rangle)).$$

The braced term is what us allows to access the linear structure of the vectors in our analysis. **Comparison to** *M***-estimators.** The works [10, 20] consider regression and sparsification for what they call general *M*-estimators. Essentially, this corresponds to the special case of our framework where all the loss functions are the same:  $f_1 = \cdots = f_m = M$ , and one assumes M(0) = 0, monotonicity, and upper and lower growth lower bounds. They additionally assume that M is p-subadditive (for p = 1/2) in the sense that  $M(x + y)^p \le M(x)^p + M(y)^p$ , which is a stronger condition than the auto-Lipschitz property a for  $h = f_i^{1/2}$ .

Under this stronger set of assumptions, the authors of [20] achieve approximations with sparsity  $\tilde{O}(n^{\max\{2,p/2+1\}})$ , which is a factor n larger than what one might hope for. In the regime  $p \le 2$  of possible near-linear-sized sparsifiers, we close this gap: Theorem 1.1 gives sparsity  $\tilde{O}(n)$ .

1.3.1 Discussion of the  $s_{\max}/s_{\min}$  Dependence. Note that Theorem 1.1 only achieves an approximation for  $s_{\min} \leq F(x) \leq s_{\max}$ , and there is a logarithmic dependence on  $s_{\max}/s_{\min}$  in the sparsity bound. Intuitively, some dependence on  $s_{\max}/s_{\min}$  is necessary in the generality of Theorem 1.1 because nothing in our assumptions precludes the functions  $f_i$  from behaving nearly independently on different scales (at least if the scales are sufficiently well separated).

In the case that each of the functions  $f_1,\ldots,f_m$  is p-homogeneous, in the sense that  $f_i(\lambda z)=|\lambda|^pf_i(z)$ , then F and the sparsifier  $\tilde{F}$  are both p-homogeneous, and therefore the guarantee  $|F(x)-\tilde{F}(x)|\leqslant \varepsilon$  for F(x)=1 already suffices to obtain  $|F(x)-\tilde{F}(x)|\leqslant \varepsilon F(x)$  for all  $x\in\mathbb{R}^n$ , meaning there is no scale dependence.

More generally, for F satisfying the hypotheses of Theorem 1.1, the growth assumptions on  $f_1, \ldots, f_m$  allow one to obtain weak guarantees even for  $F(x) \notin [s_{\min}, s_{\max}]$ . For tamer functions with only a constant number of different scaling regimes, this allows one to avoid the  $s_{\max}/s_{\min}$  dependence by applying such scaling arguments and a simple reduction. For the sake of concreteness, we demonstrate this for the Huber loss (the  $\gamma_1$  function as in (1.4)). A similar argument applies for all the  $\gamma_P$  functionals.

Lemma 1.5. Consider  $a_1, \ldots, a_m \in \mathbb{R}^n$  for  $m \ge 2$ , and  $1/m < \varepsilon < 1$ . Denote

$$F(x) := w_1 \gamma_1(\langle a_1, x \rangle) + \dots + w_m \gamma_m(\langle a_m, x \rangle)$$
  
$$\tilde{F}(x) := \tilde{w}_1 \gamma_1(\langle a_1, x \rangle) + \dots + \tilde{w}_m \gamma_m(\langle a_m, x \rangle)$$

for some nonnegative weights  $w, \tilde{w} \in \mathbb{R}^m_+$ . Suppose that

$$|F(x) - \tilde{F}(x)| \le \varepsilon F(x)$$
 for  $x \in \mathbb{R}^n$  such that  $w_{\min} \le F(x) \le 4m^2 w_{\max}$ ,

where  $w_{max} := \max(\max(w), \max(\tilde{w}))$  and  $w_{min} := \min(w)$ . Then  $\tilde{F}$  is a  $2\varepsilon$ -approximation to F.

Combining this with an analysis of the weights produced by our construction and the guarantee of Theorem 1.1 yields the following consequence. The proof of Lemma 1.5 and the next result are presented in the full version.

COROLLARY 1.6. For every  $\varepsilon > 0$ , the function

$$F(x) := \gamma_1(\langle a_1, x \rangle) + \cdots + \gamma_1(\langle a_m, x \rangle)$$

admits an s-sparse  $\varepsilon$ -approximation for

$$s \leq \frac{n}{\varepsilon^2} (\log m) \left( \log \left( \frac{n}{\varepsilon} \log m \right) \right)^3.$$

Note that our sparsity bound has an m dependence, as opposed to the classical cases of  $\ell_p$  sparsification, where sparsity bounds depend only on n and  $\varepsilon$ . However, some m dependence is not surprising, as [20, §4.5] present vectors  $a_1, \ldots, a_m \in \mathbb{R}^n$  for which the sum of the sensitivities (see (1.10)) can grow doubly-logarithmically with m:

$$\sum_{i=1}^m \max_{0 \neq x \in \mathbb{R}^n} \frac{\gamma_1(\langle a_i, x \rangle)}{F(x)} \gtrsim n \log \log \frac{m}{n}.$$

[20] also shows that  $\sum_{i=1}^{m} \max_{0 \neq x \in \mathbb{R}^n} \frac{\gamma_p(\langle a_i, x \rangle)}{F(x)} \gtrsim n \log \frac{m}{n}$  is possible for  $p \in [0, 1)$ .

# 1.4 Importance Sampling and Multiscale Weights

Given  $F(x) = f_1(\langle a_1, x \rangle) + \cdots + f_m(\langle a_m, x \rangle)$ , our approach to sparsification is via importance sampling. Given a probability vector  $\rho \in \mathbb{R}^m$  with  $\rho_1, \ldots, \rho_m > 0$  and  $\rho_1 + \cdots + \rho_m = 1$ , we sample  $M \ge 1$  coordinates  $\nu_1, \ldots, \nu_M$  i.i.d. from  $\rho$ , and define our potential approximator by

$$\tilde{F}(x) := \frac{1}{M} \sum_{i=1}^{M} \frac{f_{\nu_j}(\langle a_{\nu_j}, x \rangle)}{\rho_{\nu_j}} \,.$$

One can easily check that this gives an unbiased estimator for every  $x \in \mathbb{R}^n$ , i.e.,  $\mathbb{E}[\tilde{F}(x)] = F(x)$ .

Since we want an approximation guarantee to hold simultaneously for many  $x \in \mathbb{R}^n$ , it is natural to analyze expressions of the form

$$\mathbb{E}\max_{F(x)\leqslant s}\left|F(x)-\tilde{F}(x)\right|.$$

Analysis of this expression involves the size of discretizations of the set  $B_F(s) := \{x \in \mathbb{R}^n : F(x) \leq s\}$  at various granularities, as explained in Section 1.6.2. The key consideration (via Dudley's entropy inequality, Lemma 1.13) is how well  $B_F(s)$  can be covered by cells on which we have uniform control on how much the terms  $f_i(\langle a_i, x \rangle)/\rho_i$  vary within each cell.

**The**  $\ell_2$  **case.** Let's consider the case  $f_i(z) = |z|^2$  so that  $F(x) = |\langle a_1, x \rangle|^2 + \dots + |\langle a_m, x \rangle|^2$ . Here,  $B_F(s) = \{x \in \mathbb{R}^n : ||Ax||_2^2 \le s\}$ , where A is the matrix with  $a_1, \dots, a_m$  as rows.

A cell at scale  $2^j$  looks like

$$K_j := \left\{ x \in \mathbb{R}^n : \max_{i \in [m]} \frac{|\langle a_i, x \rangle|^2}{\rho_i} \le 2^j \right\},$$

and the pertinent question is how many translates of  $K_j$  it takes to cover  $B_F(s)$ . In the  $\ell_2$  case, this is the well-studied problem of covering Euclidean balls by  $\ell_\infty$  balls.

If  $N_j$  denotes the minimum number of such cells required, then the dual-Sudakov inequality (see Lemma 1.12) tells us that

$$\log N_j \lesssim \frac{s}{2^j} \log(m) \max_{i \in [m]} \frac{\|(A^T A)^{-1/2} a_i\|_2^2}{\rho_i}.$$

Choosing  $\rho_i := \frac{1}{n} ||(A^{\top}A)^{-1/2}a_i||_2^2$ , i.e., normalized leverage scores, yields uniform control on the size of the coverings:

$$\log N_j \lesssim \frac{s}{2i} n \log m$$
.

**The**  $\ell_p$  case,  $1 \le p < 2$ . Consider the case  $f_i(z) = |z|^p$  so that  $F(x) = ||Ax||_p^p$ . A cell at scale  $2^j$  now looks like

$$\mathsf{K}_j \coloneqq \left\{ x \in \mathbb{R}^n : \max_{i \in [m]} \frac{|\langle a_i, x \rangle|^p}{\rho_i} \leq 2^j \right\},$$

To cover  $B_F(s)$  by translates of  $K_j$ , we again employ Euclidean balls, and use  $\ell_p$  Lewis weights to relate the  $\ell_p$  structure to an  $\ell_2$ 

A classical result of Lewis [19] (see also [4, 12]) establishes that there are nonnegative weights  $w_1, \ldots, w_m \ge 0$  such that if W = $\operatorname{diag}(w_1,\ldots,w_m)$  and  $U := (A^{\top}WA)^{1/2}$ , then

$$w_i = \frac{\|U^{-1}a_i\|_2^p}{\|U^{-1}a_i\|_2^2} = \frac{f_i(\|U^{-1}a_i\|_2)}{\|U^{-1}a_i\|_2^2}.$$
 (1.5)

Assuming that A has full rank, a straightforward calculation gives  $\sum_{i=1}^{m} w_i ||U^{-1}a_i||_2^2 = \operatorname{tr}(U^2 U^{-2}) = n.$ 

Therefore, we can choose  $\rho_i := \frac{1}{n} w_i ||U^{-1} a_i||_2^2$  for i = 1, ..., m,

$$\mathsf{K}_j \coloneqq \left\{ x \in \mathbb{R}^n : \max_{i \in [m]} |\langle a_i, x \rangle|^p \leqslant \frac{2^j}{n} w_i \|U^{-1} a_i\|_2^2 \right\}.$$

(Note that the values  $\{w_i\|U^{-1}a_i\|_2^2: i=1,\ldots,m\}$  are typically referred to as the " $\ell_p$  Lewis weights".)

If we are trying to use  $\ell_2$ - $\ell_\infty$  covering bounds, we face an immediate problem: Unlike in the  $\ell_2$  case, we don't have prior control on  $||Ux||_2$  for  $x \in B_F(s)$ . One can obtain an initial bound using the structure of  $U = (A^T W A)^{1/2}$ :

$$||Ux||_{2}^{2} = \sum_{i=1}^{m} w_{i} \langle a_{i}, x \rangle^{2} \stackrel{(1.5)}{=} \sum_{i=1}^{m} ||U^{-1}a_{i}||_{2}^{p-2} \langle a_{i}, x \rangle^{2}$$

$$= \sum_{i=1}^{m} \left( \frac{|\langle a_{i}, x \rangle|}{||U^{-1}a_{i}||_{2}} \right)^{2-p} |\langle a_{i}, x \rangle|^{p} \quad (1.6)$$

$$\leq ||Ux||_{2}^{2-p} \sum_{i=1}^{m} |\langle a_{i}, x \rangle|^{p}, \quad (1.7)$$

where the last inequality is Cauchy-Schwarz:  $|\langle a_i, x \rangle| =$  $|\langle U^{-1}a_i, Ux\rangle| \le ||U^{-1}a_i||_2 ||Ux||_2$ . This gives the bound  $||Ux||_2 \le$  $||Ax||_p \leqslant s^{1/p}$  for  $x \in B_F(s)$ .

Problematically, this uniform  $\ell_2$  bound is too weak, but there is a straightforward solution: Suppose we cover  $B_F(s)$  by translates of  $K_{i_0}$ . This gives an  $\ell_{\infty}$  bound on the elements of each cell, meaning that we can apply (1.7) and obtain a better upper bound on  $||Ux||_2$ for  $x \in K_{j_0}$ . Thus to cover  $B_F(s)$  by translates of  $K_j$  with  $j < j_0$ , we will cover first by translates of  $K_{i_0}$ , then cover each translate  $(x + K_{j_0}) \cap B_F(s)$  by translates of  $K_{j_0-1}$ , and so on.

The standard approach in this setting (see [4] and [17, §15.19]) is to instead use interpolation inequalities and duality of covering numbers for a cleaner analytic version of such an iterated covering bound. However, the iterative covering argument can be adapted to the non-homogeneous setting, as we discuss next.

Generalized linear models. When we move to more general loss functions  $f_i: \mathbb{R} \to \mathbb{R}$ , we lose the homogeneity property  $f_i(\lambda x) = \lambda^p f_i(x), \lambda > 0$  that holds for  $\ell_p$  losses. Because of this, we need to replace the single Euclidean structure present in (1.5) (given by the linear operator U) with a family of structures, one for every relevant scale.

**Definition 1.7** (Approximate weights). Fix  $a_1, \ldots, a_m \in \mathbb{R}^n$  and loss functions  $f_1, \ldots, f_m : \mathbb{R} \to \mathbb{R}_+$ . We say that a vector  $w \in \mathbb{R}_+^m$ is an  $\alpha$ -approximate weight at scale s if

$$\frac{s}{\alpha} \le \frac{f_i(\|M_w^{-1/2}a_i\|_2)}{w_i\|M_w^{-1/2}a_i\|_2^2} \le \alpha s, \quad i = 1, \dots, m,$$
 (1.8)

where 
$$M_w := \sum_{j=1}^m w_j a_j a_j^{\mathsf{T}}$$
. (1.9)

To motivate this definition, let us define scale-specific sentivities:

$$\xi_i(s) := \max \left\{ \frac{f_i(\langle a_i, x \rangle)}{F(x)} : x \in \mathbb{R}^n, F(x) \in [s/2, s] \right\},$$
 (1.10)  
$$i = 1, \dots, m.$$
 (1.11)

As shown in Corollary 2.3, if the functions  $\{f_i\}$  are lower  $\theta$ homogeneous, upper 2-homogeneous, and O(1)-symmetric (in the sense of c), then an  $\alpha$ -approximate weight at scale s allows us to upper bound sensitivies by leverage scores:

$$\xi_i(s) \lesssim \sigma_i(W^{1/2}A), \qquad (1.12)$$

where  $W = \text{diag}(w_1, \dots, w_m)$ , and the implicit constant depends on  $\alpha$  and the homogeneity parameters. Here,  $\sigma_i(V)$  denotes the *i*th leverage score of a matrix V with rows  $v_1, \ldots, v_m$ :

$$\sigma_i(V) := \langle v_i, (V^\top V)^+ v_i \rangle, \tag{1.13}$$

where  $(V^{\top}V)^{+}$  denotes the Moore-Penrose pseudoinverse. Notably, one always has  $\sigma_1(V) + \cdots + \sigma_m(V) = \text{rank}(V)$ , and therefore (1.12) gives an upper bound  $\xi_1(s) + \cdots + \xi_m(s) \leq n$ .

In order to generalize the iterated covering argument for  $\ell_p$  losses, we need there to be a relationship between weights at different

**Definition 1.8** (Weight schemes). Let  $\mathcal{J} \subseteq \mathbb{Z}$  be a contiguous interval. A family  $\{w^{(j)} \in \mathbb{R}^m_+ : j \in \mathcal{J}\}$  is an  $\alpha$ -approximate weight scheme if each  $w^{(j)}$  is an  $\alpha$ -approximate weight at scale  $2^{j}$ and, furthermore, for every pair  $j, j + 1 \in \mathcal{J}$  and  $i \in \{1, ..., m\}$ ,

$$w_i^{(j+1)} \leqslant \alpha w_i^{(j)} \,. \tag{1.14}$$

Given a weight scheme, we choose sampling probabilities

$$\rho_i \propto \max_{i \in \mathcal{T}} w_i^{(j)} \| M_{w^{(j)}}^{-1/2} a_i \|_2^2 = \max_{i \in \mathcal{T}} \sigma_i(W_j^{1/2} A), \quad i = 1, \dots, m,$$

where  $W_j = \operatorname{diag}(w_1^{(j)}, \dots, w_m^{(j)})$ . In our setting,  $|\mathcal{J}| \leq$  $O(\log(ms_{\max}/s_{\min}))$ , which results in the sparsity increasing by a corresponding factor.

In Section 2, we establish the existence of approximate weight schemes for general families of loss functions satisfying certain growth bounds, along with efficient algorithms to compute the corresponding weights.

# Regression via Iterative Refinement

Previous works have observed that combining iterative refinement with sparsification of  $\gamma_p$ -functions (recall (1.4)) leads to improved algorithms for  $\ell_p$ -regression [1, 2, 13]. For the benefit of the reader, we give a description of these ideas in somewhat more generality.

Recall that our goal is to find a point  $x \in \mathbb{R}^n$  that computes an approximate minimizer of  $F(x) := \sum_{i=1}^{m} f_i(\langle a_i, x \rangle - b_i)$ , up to high accuracy. For now, we assume that F is a differentiable convex function and denote  $F_* := \inf_{x \in \mathbb{R}^n} F(x)$ . Later, we will introduce additional conditions that allow for iterative refinement to succeed.

Broadly, iterative refinement minimizes F(x) be repeatedly solving sub-problems, each of which make multiplicative progress in reducing the error of the current solution. Given a current point  $x_0$ , prior works on iterative refinement define a local approximation of F suitably symmetrized and centered around  $x_0$  such that approximately minimizing this local approximation yields the desired decrease in function error. One way to derive such local approximations is through Bregman divergences, which give a natural way of recentering convex functions.

**Definition 1.9** (The *F*-divergence). For  $x, y \in \mathbb{R}^n$ , use  $T_x^F(y) := F(x) + \nabla F(x)^\top (y-x)$  to denote the first order Taylor approximation of *F* at *x*, and define the *F*-induced Bregman divergence by  $D_x^F(y) := F(y) - T_x^F(y)$ .

Note that for convex F, the function  $D_x^F(y)$  is convex and minimized at x. Consequently, given a point  $x_0$ , the function  $D_{x_0}^F(y)$  is a natural function induced by F and minimized at  $x_0$ . Note that minimizing F(x) is the same as minimizing  $\langle \nabla F(x_0), x - x_0 \rangle + D_{x_0}^F(x)$  which in turn is the same as minimizing  $\langle \nabla F(x), \Delta \rangle + D_{x_0}^F(x_0 + \Delta)$  over  $\Delta$  and adding the minimizer to  $x_0$ .

Iterative refinement strategies approximately minimize  $\langle \nabla F(x), \Delta \rangle + r(\Delta)$ , where  $r(\Delta)$  is a suitable approximation of  $D^F_{x_0}(x_0 + \Delta)$ . One step of refinement moves to  $x_1 := x_0 + \eta \tilde{\Delta}$ , where  $\eta$  is a suitably chosen step-size and  $\tilde{\Delta}$  is the approximate minimizer. As motivation for our approach, here we consider the scheme suggested by prior work, where  $r(\Delta)$  is a sparsification of a simple approximation to the divergence.

Informally, one can show that if the square root of the Bregman divergence of each  $f_i$  is L-auto-Lipschitz a and lower  $\theta$ -homogeneous for some  $\theta > 1$  b, then one step of sparsification/refinement decreases the error in the objective value multiplicatively by an absolute constant. Interestingly, auto-Lipschitzness is only required for sparsification and not for refinement. However, we critically need the Bregman divergence to be lower  $\theta$ -homogeneous for  $\theta > 1$  for iterative refinement, while our sparsification results (Theorem 1.1) only require  $\theta > 0$ .

LEMMA 1.10 (REFINEMENT LEMMA). Suppose  $r: \mathbb{R}^n \to \mathbb{R}_+$  is lower  $\theta$ -homogeneous with constant c < 1 for  $\theta > 1$ , and for  $x_0 \in \mathbb{R}^n$  and all  $\Delta \in \mathbb{R}^n$  and  $\eta \in [0, 1]$ ,

$$r(\Delta) \leq D^F_{x_0}(x_0 + \Delta) \leq \alpha \, r(\Delta)$$

where  $\alpha \ge 1$  is fixed. Then  $\inf_{\Delta \in \mathbb{R}^n} \left\{ T_{x_0}^F(x_0 + \hat{\Delta}) + r(\hat{\Delta}) \right\} \le F_*$  and if  $\hat{\Delta} \in \mathbb{R}^n$  satisfies

$$T_{x_0}^F(x_0 + \hat{\Delta}) + r(\hat{\Delta}) \le F_*,$$
 (1.15)

then

$$F(x_0 + \hat{\eta}\hat{\Delta}) - F_* \le (1 - \hat{\eta}) (f(x_0) - F_*), \tag{1.16}$$

where 
$$\hat{\eta} := (\alpha/c)^{-1/(q-1)}$$
. (1.17)

PROOF. First note that, for all  $\Delta \in \mathbb{R}^n$ ,

$$F(x_0 + \Delta) = T_{x_0}^F(x_0 + \Delta) + D_{x_0}^F(x_0 + \Delta). \tag{1.18}$$

Since,  $D_{x_0}^F(x_0 + \Delta) \ge r(\Delta)$  this implies the desired bound

$$\inf_{\Delta \in \mathbb{R}^n} \left\{ T_{x_0}^F(x_0 + \Delta) + r(\Delta) \right\} \leqslant \inf_{\Delta \in \mathbb{R}^n} F(x_0 + \Delta) = F_*.$$

Next, note that for all  $\eta \in [0, 1]$  and  $\Delta \in \mathbb{R}^n$ ,

$$T_{x_0}^F(x_0 + \eta \Delta) = F(x_0) + \eta \nabla F(x_0)^{\top} \Delta$$
 (1.19)

$$= (1 - \eta)F(x_0) + \eta T_{x_0}^F(x_0 + \Delta)$$
 (1.20)

$$D_{x_0}^F(x_0 + \eta \Delta) \le \alpha r(\eta \Delta) \le \alpha/c \cdot \eta^{\theta} r(\Delta)$$
 (1.21)

Suppose that  $\hat{\Delta} \in \mathbb{R}^n$  satisfies (1.15). Then plugging (1.19) and (1.21) into (1.18) with  $\Delta = \hat{\eta} \hat{\Delta}$  yields

$$\begin{split} F(x_0 + \hat{\eta} \hat{\Delta}) & \leq (1 - \hat{\eta}) F(x_0) + \hat{\eta} \, T_{x_0}^F(x_0 + \hat{\Delta}) + \alpha/c \cdot \hat{\eta}^{\theta} r(\hat{\Delta}) \\ & = (1 - \hat{\eta}) F(x_0) + \hat{\eta} \left( T_{x_0}^F(x_0 + \hat{\Delta}) + r(\hat{\Delta}) \right) \\ & \leq (1 - \hat{\eta}) F(x_0) + \hat{\eta} F_* \end{split}$$

where we used that  $\alpha/c \cdot \hat{\eta}^{\theta-1} = 1$  and  $\hat{\eta} \in [0, 1]$ . Rearranging yields (1.16).

To apply Lemma 1.10 to ERM for general linear models, note that

$$D_{x_0}^F(x_0 + \Delta) = \sum_{i=1}^m D_{\langle a_i, x_0 \rangle - b_i}^{f_i} (\langle a_i, \Delta \rangle).$$

Now, if the square root of each divergence  $D_z^{f_i}$  is L-auto-Lipschitz and lower  $\theta$ -homogeneous, then Theorem 1.1 gives weights  $w \in \mathbb{R}_+^m$  with sparsity  $\tilde{O}(n)$  such that

$$0.9 \cdot D_{x_0}^F(x_0 + \Delta) \leq r(\Delta) \leq D_{x_0}^F(x_0 + \Delta)$$
 where 
$$r(\Delta) \coloneqq \sum_{i=1}^m w_i D_{\langle a_i, x_0 \rangle - b_i}^f(\langle a_i, \Delta \rangle).$$

Thus, Lemma 1.10 applies, and we can decrease the objective value error by a multiplicative factor by minimizing  $r(\Delta)$ . Since  $r(\Delta)$  only has  $\tilde{O}(n)$  nonzero terms, one can apply previous solvers for  $m = \tilde{O}(n)$  [2, 18] to obtain the desired runtimes.

In the full version's Section on  $\ell_p$  regression, we verify that the f-divergence of  $f(z) = |z|^p$  is the  $\gamma_p$  function, which is lower p-homogeneous, and has an auto-Lipschitz square root. This yields our algorithm for  $\ell_p$  regression (Theorem 1.2). A formal version of the argument is presented in the full version. The main difference is that some technical work is needed because Theorem 1.1 only provides sparsification for a range of inputs  $\{x \in \mathbb{R}^n : s_{\min} \le F(x) \le s_{\max}\}$ . Moreover, our algorithm uses an approximate oracle for GLMs with the functions  $\{f_i\}$  (rather than with the approximate divergence  $r(\cdot)$ ), and we show that the oracle does not need to be solved to high accuracy to make sufficient progress.

### 1.6 Preliminaries

Throughout the paper, we denote  $[n] := \{1, 2, \dots, n\}$ . We use the notation  $a \leq b$  to denote that there is a universal constant C such that  $a \leq Cb$ , and  $a \leq_L b$  to denote that C may depend on L. We use  $a \times b$  to denote the conjunction of  $a \leq b$  and  $b \leq a$ , and  $a \times_L b$  analogously. We also denote  $\mathbb{R}_+ := \{x \in \mathbb{R} : x \geq 0\}$  and  $\mathbb{R}_{++} := \{x \in \mathbb{R} : x > 0\}$ .

For simplicity of presentation, we assume that the vectors  $a_1, \ldots, a_m \in \mathbb{R}^n$  in (1.1) span  $\mathbb{R}^n$  and all are nonzero. In particular, this means that the matrix A with rows  $a_1, \ldots, a_m$  has rank n and  $A^TWA$  is invertible for any diagonal matrix W with strictly positive entries on the diagonal.

1.6.1 Covering numbers and chaining. Consider a metric space (T,d). For  $x \in T$  and r > 0, define the ball  $B(x,r) := \{y \in T : d(x,y) \le r\}$ .

**Definition 1.11** (Covering numbers). For a radius r > 0, we define the *covering number*  $\mathcal{N}(T,d,r)$  as the smallest number of balls of radius r (in the distance d) that are required to cover T. For  $S,S' \subseteq \mathbb{R}^n$ , we overload notation and use  $\mathcal{N}(S,S')$  to denote the smallest number of translates of S' needed to cover S.

We require the following "dual Sudakov inequality" (see [21] and [17, (3.15)]) which gives bounds for covering the Euclidean ball using balls in an arbitrary norm.

LEMMA 1.12 (DUAL SUDAKOV INEQUALITY). Let  $B_2^n$  denote the unit ball in n dimensions, and  $\|\cdot\|_X$  an arbitrary norm on  $\mathbb{R}^n$ . If g is a standard n-dimensional Gaussian, then

$$\sqrt{\log \mathcal{N}(B_2^n, B_X)} \lesssim \mathbb{E} \| \boldsymbol{g} \|_X,$$

where  $B_X := \{ y \in \mathbb{R}^n : ||y||_X \le 1 \}.$ 

We recall Talagrand's generic chaining functional [27, Def. 2.2.19]:

$$\gamma_2(T,d) := \inf_{\{\mathcal{A}_h\}} \sup_{x \in T} \sum_{h=0}^{\infty} 2^{h/2} \operatorname{diam}(\mathcal{A}_h(x), d), \qquad (1.22)$$

where the infimum runs over all sequences  $\{\mathcal{A}_h : h \ge 0\}$  of partitions of T satisfying  $|\mathcal{A}_h| \le 2^{2^h}$  for each  $h \ge 0$ . Note that we use the notation  $\mathcal{A}_h(x)$  for the unique set of  $\mathcal{A}_h$  that contains x.

The chaining functional is used to control the maximum of subgaussian processes (see, e.g., the discussion in [14, §2.2] where it is applied precisely in the setting of sparsification). Our use of the functional occurs only in our application of our abstract sparsification statement (stated in the full version), and in this paper we will only require the following classical upper bound. (See, eg., [27, Prop 2.2.10].)

Lemma 1.13 (Dudley's entropy bound). For any metric space (T,d), it holds that

$$\gamma_2(T,d) \lesssim \sum_{j \in \mathbb{Z}} 2^j \sqrt{\log \mathcal{N}(T,d,2^j)}$$
.

The interested reader will note that this is (up to constants) precisely the upper bound one obtains by choosing  $\mathcal{A}_h$  as a uniform discretization of (T,d), i.e., to minimize  $\sup\{\operatorname{diam}(\mathcal{A}_h(x),d):x\in T\}$  over all partitions satisfying  $|\mathcal{A}_h|\leqslant 2^{2^h}$ .

1.6.2 Sparsification via subgaussian processes. We discuss sparsification via subgaussian processes. Consider

$$\varphi_1, \varphi_2, \ldots, \varphi_m : \mathbb{R}^n \to \mathbb{R},$$

and define

$$F(x) := \sum_{j=1}^{m} \varphi_j(x) .$$

Given a strictly positive probability vector  $\rho \in \mathbb{R}^m_{++}$ , and an integer  $s \ge 1$  and  $\nu = (\nu_1, \dots, \nu_s) \in [m]^s$ , define the distance

$$d_{\rho,\nu}(x,y) := \left(\sum_{j=1}^{s} \left(\frac{\varphi_{\nu_j}(x) - \varphi_{\nu_j}(y)}{\rho_{\nu_j} s}\right)^2\right)^{1/2}.$$
 (1.23)

and the function  $\tilde{F}_{\rho,\nu}:\mathbb{R}^n\to\mathbb{R}$ 

$$\tilde{F}_{\rho,\nu}(x) := \frac{1}{s} \sum_{i=1}^s \frac{\varphi_{\nu_j}(x)}{\rho_{\nu_j}}.$$

#### 2 MULTISCALE IMPORTANCE SCORES

Recall the definitions of approximate weights (Definition 1.7) and weight schemes (Definition 1.8). In the present section, we prove the following two results.

THEOREM 2.1. Suppose that  $f_1, \ldots, f_m : \mathbb{R} \to \mathbb{R}_+$  are lower  $\theta$ -homogeneous and upper u-homogeneous with  $u > \theta > 0$  and uniform constants c, C > 0. Then there is some  $\alpha = \alpha(\theta, c, u, C) > 1$  such that for every choice of vectors  $a_1, \ldots, a_m \in \mathbb{R}^n$  and s > 0, there is an  $\alpha$ -approximate weight at scale s.

This is proved in Section 2.2 by considering critical points of the functional  $U \mapsto \det(U)$  subject to the constraint  $G(U) \leq s$ , where  $G(U) \coloneqq f_1(\|Ua_1\|_2) + \cdots + f_m(\|Ua_m\|_2)$ , which can be seen as a generalization of Lewis' original method.

**Single-scale sensitivities.** Let us now observe that, in the case  $u \le 2$ , Theorem 2.1 allows us to bound sensitivities (recall (1.10)). The next lemma is a generalization of (1.7).

LEMMA 2.2. Suppose  $f_1, \ldots, f_m : \mathbb{R} \to \mathbb{R}_+$  satisfy the assumptions of Theorem 2.1 with  $u \leq 2$ , and they are additionally K-symmetric in the sense of c. If  $w \in \mathbb{R}_+^m$  is an  $\alpha$ -approximate weight at scale s, then for any  $x \in \mathbb{R}^n$ , it holds that

$$||M_w^{1/2}x||_2^{\theta} \le \max\left(1, \alpha \frac{CK}{c} \frac{F(x)}{s}\right).$$

PROOF. We may clearly assume that  $||M_w^{1/2}x||_2 \ge 1$ . Then using the Cauchy-Schwarz inequality

$$|\langle a_i, x \rangle| \le ||M_w^{-1/2} a_i||_2 ||M_w^{1/2} x||_2$$

together with the upper quadratic growth assumption gives

$$\frac{f_i(|\langle a_i, x \rangle|)}{|\langle a_i, x \rangle|^2} \ge \frac{1}{C} \frac{f_i(||M_w^{1/2} x||_2 ||M_w^{-1/2} a_i||_2)}{||M_w^{1/2} x||_2^2 ||M_w^{-1/2} a_i||_2^2}$$
(2.1)

$$\geq \frac{c}{C} \frac{\|M_{\mathbf{w}}^{1/2} x\|_{2}^{\theta} f_{i}(\|M_{\mathbf{w}}^{-1/2} a_{i}\|_{2})}{\|M_{\mathbf{w}}^{1/2} x\|_{2}^{2} \|M_{\mathbf{w}}^{-1/2} a_{i}\|_{2}^{2}}, \qquad (2.2)$$

where the last inequality uses the lower growth assumption. Using  $M_w = \sum_{i=1}^m w_i a_i a_i^{\mathsf{T}}$ , we can bound

$$||M_{w}^{1/2}x||_{2}^{2} = \sum_{i=1}^{m} w_{i} \langle a_{i}, x \rangle^{2}$$

$$\stackrel{(2.1)}{\leq} \frac{C}{c} ||M_{w}^{1/2}x||_{2}^{2-\theta} \sum_{i=1}^{m} w_{i} f_{i}(|\langle a_{i}, x \rangle|) \frac{||M_{w}^{-1/2}a_{i}||_{2}^{2}}{f_{i}(||M_{w}^{-1/2}a_{i}||_{2})}$$

$$\leq \alpha \frac{CK}{c} \frac{1}{s} \|M_w^{1/2} x\|_2^{2-\theta} \sum_{i=1}^m f_i(\langle a_i, x \rangle),$$

where the last inequality uses the defining property of an  $\alpha$ -approximate weight at scale s, along with the assumption of K-symmetry:  $f_i(|\langle a_i, x \rangle|) \leq K f_i(\langle a_i, x \rangle)$  for all i = 1, ..., m.

COROLLARY 2.3 (SENSITIVITY UPPER BOUND). Under the assumptions of Lemma 2.2, it holds that

$$\xi_1(s) + \cdots + \xi_m(s) \lesssim n$$
,

where the implicit constant depends on the parameters  $\alpha, \theta, C, c, K$ .

PROOF. From Lemma 2.2, if  $F(x) \le s$ , then  $||M_w^{1/2}x||_2 \le (\alpha CK/c)^{1/\theta}$ . By Cauchy-Schwarz, this gives

$$|\langle a_i, x \rangle| \leq (\alpha C K/c)^{1/\theta} ||M_w^{-1/2} a_i||_2$$

and therefore

$$\begin{split} & f_i(\langle a_i, x \rangle) \leqslant K f_i(|\langle a_i, x \rangle|) \\ & \leqslant \frac{K}{c} f_i \left( (\alpha C K/c)^{1/\theta} \| M_w^{-1/2} a_i \|_2 \right) \\ & \leqslant \frac{CK}{c} \left( \frac{\alpha C K}{c} \right)^{2/\theta} f_i(\| M_w^{-1/2} a_i \|_2) \\ & \leqslant \frac{\alpha C K}{c} \left( \frac{\alpha C K}{c} \right)^{2/\theta} s \cdot w_i \| M_w^{-1/2} a_i \|_2^2, \end{split}$$

where the first inequality uses K-symmetry, the second uses 1/c-monotonicity (which holds by Lemma 1.4), the third inequality uses upper homogeneity, and the last inequality uses that w is an  $\alpha$ -approximate weight at scale s. Finally, one notes that  $\sum_{i=1}^m w_i \|M_w^{-1/2}a_i\|_2^2 = \operatorname{tr}\left(M_w^{-1}M_w\right) = n$ , completing the proof.  $\square$ 

We are only able to establish the existence of entire weight schemes (where the weights at adjacent scales are related) for u < 4, which suffices our applications, as  $u \le 2$  is a requirement for Theorem 1.1. The following theorem is proved in Section 2.1, based on the contractive iteration method introduced by Cohen and Peng [12].

THEOREM 2.4. Suppose that  $f_1, \ldots, f_m : \mathbb{R} \to \mathbb{R}_+$  are lower  $\theta$ -homogeneous and upper u-homogeneous with  $4 > u > \theta > 0$  and uniform constants c, C > 0. Then there is some  $\alpha = \alpha(u, C, \theta, c)$  such that for every choice of vectors  $a_1, \ldots, a_m \in \mathbb{R}^n$ , there is an  $\alpha$ -approximate weight scheme  $\{w_i^{(j)}: j \in \mathbb{Z}\}$ .

In the next section, we show how to compute an approximate weight scheme  $\{w_i^{(j)}: j\in \mathcal{J}\}$  using  $\tilde{O}(|\mathcal{J}|)$  computations of leverage scores  $(\sigma_1(V),\ldots,\sigma_m(V))$  for matrices of the form  $V=A^\top WA$ .

#### 2.1 Contractive Algorithm

For a weight  $w \in \mathbb{R}^m_+$  and  $i \in \{1, ..., m\}$ , define

$$\tau_i(w) := \frac{\sigma_i(W^{1/2}A)}{w_i} = \langle a_i, (A^\top WA)^{-1} a_i \rangle,$$

$$W := \operatorname{diag}(w_i, w_i)$$

and denote  $\tau(w) := (\tau_1(w), \dots, \tau_m(w)).$ 

Fix a scale parameter s>0 and define the iteration  $\varphi_s:\mathbb{R}^m_+\to\mathbb{R}^m_+$  by

$$(\varphi_s(w))_i := \frac{1}{s} \frac{f_i(\sqrt{\tau_i(w)})}{\tau_i(w)}. \tag{2.3}$$

Write  $\varphi^k := \varphi \circ \cdots \circ \varphi$  for the k-fold composition of  $\varphi$ . In this case where  $f_i(z) = |z|^p$  and  $1 \le p \le 2$ , it is known, for s = 1, starting from any  $w_0 \in \mathbb{R}_+^m$ , the sequence  $\{\varphi_1^k(w_0) : k \ge 1\}$  converges to the unique fixed point of  $\varphi$ , which are the corresponding  $\ell_p$  Lewis weights (1.5).

Define now a metric d on  $\mathbb{R}^m_+$  by

$$d(u, w) := \max \left\{ \left| \log \frac{u_i}{w_i} \right| : i = 1, \dots, m \right\}.$$

We note the following characterization.

FACT 2.5. A vector  $w \in \mathbb{R}^m_+$  is an  $\alpha$ -approximate weight at scale s if and only if

$$d(w, \varphi_s(w)) \leq \log \alpha$$
.

First, we observe that  $\tau$  is 1-Lipschitz on  $(\mathbb{R}^m_+, d)$ . In the next proof,  $\leq$  denotes the ordering of two real, symmetric matrices in the Loewner order, i.e.,  $A \leq B$  if and only if B - A is positive semi-definite.

Lemma 2.6. For any  $w, w' \in \mathbb{R}^m_+$ , it holds that

$$d(\tau(w), \tau(w')) \le d(w, w').$$

Proof. Denote

 $W = \operatorname{diag}(w), W' = \operatorname{diag}(w'), \text{ and } \alpha := \exp(d(w, w')).$ 

Then  $\alpha^{-1}W \leq W' \leq \alpha W$ , therefore

$$\alpha^{-1}A^{\mathsf{T}}WA \leq A^{\mathsf{T}}W'A \leq \alpha A^{\mathsf{T}}WA$$

and by monotonicity of the matrix inverse in the Loewner order,  $\alpha^{-1}(A^{\top}WA)^{-1} \leq (A^{\top}W'A)^{-1} \leq \alpha(A^{\top}WA)^{-1}$ . This implies  $d(\tau(w), \tau(w')) \leq \log \alpha$ , completing the proof.

PROOF OF THEOREM 2.4. Consider the map  $\psi: \mathbb{R}^m_+ \to \mathbb{R}^m_+$  whose *i*-th coordinate is defined as

$$\psi_i(x) \coloneqq \frac{f_i(\sqrt{x_i})}{x_i} \, .$$

Our assumptions on lower and upper-homogeneity give, for all  $y_i \ge x_i$ ,

$$c\left(\frac{y_i}{x_i}\right)^{\theta/2-1} \leqslant \frac{f_i(\sqrt{y_i})/y_i}{f_i(\sqrt{x_i})/x_i} \leqslant C\left(\frac{y_i}{x_i}\right)^{u/2-1},$$

yielding, for  $C_1 := \max\{C, 1/c\}$ 

$$d(\psi(x), \psi(y)) \le \max\left(\left|\frac{\theta}{2} - 1\right|, \left|\frac{u}{2} - 1\right|\right) d(x, y) + \log(C_1). \quad (2.4)$$

Fix s>0 and consider the mapping  $\varphi:\mathbb{R}^m_+\to\mathbb{R}^m_+$  defined in (2.3). Then for u<4 and  $\delta:=\max\left(\left|\frac{\theta}{2}-1\right|,\left|\frac{u}{2}-1\right|\right)<1$ , (2.4) in conjunction with Lemma 2.6, shows that

$$d(\varphi_{\mathcal{S}}(w), \varphi_{\mathcal{S}}(w')) < \delta d(w, w') + \log(C_1). \tag{2.5}$$

Applying this bound inductively, for any weight  $w \in \mathbb{R}^m_+$  and  $k \ge 1$ , we have

$$d\left(\varphi_s^k(w), \varphi_s^{k+1}(w)\right) \leq \frac{\delta^k d(\varphi_s(w), w) + \log C_1}{1 - \delta}, \qquad (2.6)$$

Now define

$$w^{(0)} \coloneqq \varphi_1^k(1,\ldots,1),$$

where  $k \ge 1$  is chosen large enough so that

$$d(w^{(0)}, \varphi_1(w^{(0)})) \leq \frac{2 \log C_1}{1 - \delta}.$$

From Fact 2.5, one sees that  $w^{(0)}$  is an  $\alpha$ -approximate weight at scale 1 for  $\alpha = C_1^{2/(1-\delta)}$ .

Define inductively, for j = 1, 2, ...,

$$\begin{split} w^{(j)} &\coloneqq \varphi_{2^j}(w^{(j-1)}) \\ w^{(-j)} &\coloneqq \varphi_{2^{-j}}(w^{(1-j)}) \,. \end{split}$$

Note that

$$d(\varphi_{2^{j}}(w^{(j)}), w^{(j)}) =$$

$$d(\varphi_{2^{j}}^{2}(w^{(j-1)}), \varphi_{2^{j}}(w^{(j-1)}))$$

$$\leq \delta d(\varphi_{2^{j}}(w^{(j-1)}), w^{(j-1)}) + \log(C_{1})$$

$$\leq \delta d(\varphi_{2^{j-1}}(w^{(j-1)}), w^{(j-1)}) + \delta \log(2) + \log(C_{1}),$$

where the last inequality uses  $\varphi_{2s}(w) = 2\varphi_s(w)$  for all  $w \in \mathbb{R}_+^m$ . Therefore, by induction,

$$d(\varphi_{2^j}(w^{(j)}), w^{(j)}) \le \frac{2\log(C_1) + \log 2}{1 - \delta}$$

for all j > 0. To see that the family of weights  $\{w^{(j)} : j \in \mathbb{Z}\}$  forms a weight scheme, note that

$$d(w^{(j)}, w^{(j-1)}) = d(\varphi_{2^{j}}(w^{(j-1)}), w^{(j-1)})$$
  
$$\leq d(\varphi_{2^{j}}(w^{(j-1)}), w^{(j-1)}) + \log 2,$$

thus  $\{w^{(j)}: j \in \mathbb{Z}\}$  is an  $\alpha$ -approximate weight scheme for  $\alpha = \frac{2\log(2C_1)}{1-\delta}$ , completing the proof.

# 2.2 A Variational Approach to Approximate Weights

In this section we show that even if each  $f_i$  is upper u-homogeneous for  $u \ge 4$ , approximate weights still exist. Our sparsification analysis relies on the existence of weight schemes, where there is a relationship between weights at different scales. But the following existence proof is instructive.

THEOREM 2.7. Suppose  $f_1, \ldots, f_m : \mathbb{R}_+ \to \mathbb{R}_+$  are lower  $\theta$ -homogeneous with constant c and upper u-homogeneous with constant c with c0. Then there is a constant c0 and c0 such that for every choice of vectors c1, ..., c2 and c3 on there is an c4-approximate weight at scale c5.

The idea behind the proof is to set up a variational problem whose critical points produce approximate weights at a given scale. As observed in [22], this analysis technique does not require convexity.

Lemma 2.8. Suppose  $g_1, \ldots, g_m : \mathbb{R}_+ \to \mathbb{R}_+$  are monotone increasing, continuously differentiable, and satisfy  $g_1(0) = \cdots = g_m(0) = 0$ . Then for every  $\beta > 0$ , there are weights  $\{w_i \ge 0 : i = 1, \ldots, m\}$  such that

$$w_i = \gamma \frac{g_i'(\|M_w^{-1/2}a_i\|_2)}{\|M_w^{-1/2}a_i\|_2} \quad i = 1, \dots, m,$$
(2.7)

$$\gamma = n \left( \sum_{i=1}^{m} g_i'(\|M_w^{-1/2} a_i\|_2) \|M_w^{-1/2} a_i\|_2 \right)^{-1},$$

$$\beta = \sum_{i=1}^{m} g_i(\|M_w^{-1/2} a_i\|_2).$$
(2.8)

PROOF. For a linear operator  $U: \mathbb{R}^n \to \mathbb{R}^n$ , define

$$G(U) := \sum_{i=1}^{m} g_i(\|Ua_i\|_2),$$

and consider the optimization

maximize 
$$\{\det(U): G(U) \le \beta\}$$
. (2.9)

Since G(0) = 0 and each  $g_i$  is monotone increasing, it holds that  $G(cI) = \beta$  for some c > 0. Therefore for any maximizer  $U^*$ , it holds that  $\det(U^*) > 0$ , i.e.,  $U^*$  is invertible, and  $G(U^*) = \beta$ .

For U invertible, we have

$$\nabla \det(U) = \det(U)U^{-\top}. \tag{2.10}$$

Let us also calculate

$$dG(U) = \sum_{i=1}^{m} g'_{i}(\|Ua_{i}\|_{2}) \frac{d\|Ua_{i}\|_{2}^{2}}{2\|Ua_{i}\|_{2}},$$

and use  $||Ua_i||_2^2 = \operatorname{tr}(U^{\top}Ua_ia_i^{\top})$  to write

$$\frac{1}{2}\mathsf{d}\|Ua_i\|_2^2 = \mathsf{tr}\left((\mathsf{d}U)^\top Ua_ia_i^\top\right),$$

so that

$$\nabla_{IJ}G(U) = UA^{\top}D_{IJ}A,$$

where  $D_U$  is the  $m \times m$  diagonal matrix with  $(D_U)_{ii} = \frac{g_i'(\|Ua_i\|_2)}{\|Ua_i\|_2}$  and  $A \in \mathbb{R}^{m \times n}$  is the matrix with rows  $a_1, \dots, a_m$ .

Combined with (2.10), we see that if U is an optimal solution to (2.9), then for some Lagrange multiplier  $\gamma > 0$ , we have

$$(U^{\top}U)^{-1} = \gamma A^{\top}D_{U}A.$$

Take  $V \coloneqq (U^\top U)^{-1}$  so that  $V = \sum_{i=1}^m w_i a_i a_i^\top$  with  $w_i \coloneqq \gamma \frac{g_i'(\|V^{-1/2}a_i\|_2)}{\|V^{-1/2}a_i\|_2}$ . To compute the value of  $\gamma$ , calculate

$$n = \operatorname{tr}(VV^{-1}) = \gamma \sum_{i=1}^{m} \frac{g_i'(\|V^{-1/2}a_i\|_2)}{\|V^{-1/2}a_i\|_2} \operatorname{tr}(V^{-1}a_ia_i^{\top})$$
$$= \gamma \sum_{i=1}^{m} g_i'(\|V^{-1/2}a_i\|_2) \|V^{-1/2}a_i\|_2.$$

The preceding lemma assumes that the functions  $g_i$  are monotone increasing. However, the functions  $f_i$  only satisfy lower homogeneity, which is a weaker condition. To prove Theorem 2.7, one can take monotone approximations of functions  $f_i$  by averaging over intervals.

Lemma 2.9. Suppose f is lower  $\theta$ -homogeneous with constant c and upper u-homogeneous with constant C for  $u \ge 1$ . Define  $K := \max(e, (2/c)^{1/\theta})$  and  $g : \mathbb{R}_+ \to \mathbb{R}_+$  by

$$g(x) := \int_{r}^{Kx} \frac{f(t)}{t} dt.$$

Then g is continuously differentiable, monotone increasing, and satisfies g(0) = 0. Moreover, for all x > 0,

$$\frac{f(x)}{x} \le g'(x) \le CK^u \frac{f(x)}{x}, \quad \forall x \ge 0$$
$$c f(x) \le g(x) \le CK^u f(x).$$

PROOF. Note that  $g'(x) = \frac{f(Kx) - f(x)}{x}$ , and  $f(Kx) \ge cK^{\theta}f(x) \ge 2f(x)$  by lower  $\theta$ -homogeneity. Thus g is monotone increasing. Moreover, we have

$$g(x) = \int_{x}^{Kx} \frac{f(t)}{t} dt \ge c f(x) \int_{x}^{Kx} \frac{1}{t} dt \ge c \cdot f(x),$$

and for  $u \ge 1$ .

$$g(x) = \int_x^{Kx} \frac{f(t)}{t} dt \le Cf(x) \int_x^{Kx} \frac{(t/x)^u}{t} dt \le CK^u f(x).$$

Proof of Theorem 2.7. Throughout the proof, we use  $\times$  in place of  $\times_{\theta,c,u,C}$ . Let K be as in Lemma 2.9 and define  $g_i(x) \coloneqq \int_x^{Kx} \frac{f_i(t)}{t} dt$ . By Lemma 2.9, we have  $g_i'(x) \times f_i(x)/x$  and  $g(x) \times f_i(x)$  for  $i=1,\ldots,m$ .

Let  $\{w_i\}$  and  $M_w$  be as in Lemma 2.8 when applied to  $g_1, \ldots, g_m$  with  $\beta = n/\lambda$ . Then,

$$\sum_{i=1}^{m} g_i'(\|M_w^{-1/2}a_i\|_2)\|M_w^{-1/2}a_i\|_2 \approx \sum_{i=1}^{m} f_i(\|M_w^{-1/2}a_i\|_2)$$
$$\approx \sum_{i=1}^{m} g_i(\|M_w^{-1/2}a_i\|_2)$$
$$= \beta = \frac{n}{\lambda}.$$

Therefore  $\gamma \approx \lambda$  (recall (2.8)), and thus (2.7) gives the desired result.

#### **ACKNOWLEDGMENTS**

Part of this work was conducted while the authors were visiting the Simons Institute for the Theory of Computing. James R. Lee is supported in part by NSF CCF-2007079 and a Simons Investigator Award. Yang P. Liu is partially supported by the Google Research Fellowship and NSF DMS-1926686. Aaron Sidford was supported in part by a Microsoft Research Faculty Fellowship, NSF CCF-1844855, NSF CCF-1955039, a PayPal research award, and a Sloan Research Fellowship.

# **REFERENCES**

- [1] Deeksha Adil, Brian Bullins, Rasmus Kyng, and Sushant Sachdeva. 2021. Almost-Linear-Time Weighted ℓ<sub>p</sub>-Norm Solvers in Slightly Dense Graphs via Sparsification. In 48th International Colloquium on Automata, Languages, and Programming, ICALP 2021, July 12-16, 2021, Glasgow, Scotland (Virtual Conference) (LIPIcs, Vol. 198), Nikhil Bansal, Emanuela Merelli, and James Worrell (Eds.). Schloss Dagstuhl Leibniz-Zentrum für Informatik, 9:1-9:15. https://doi.org/10.4230/LIPIcs.ICALP.2021.9
- [2] Deeksha Adil, Rasmus Kyng, Richard Peng, and Sushant Sachdeva. 2019. Iterative Refinement for ℓp-norm Regression. In Proceedings of the Thirtieth Annual ACM-SIAM Symposium on Discrete Algorithms, SODA 2019, San Diego, California, USA, January 6-9, 2019, Timothy M. Chan (Ed.). SIAM, 1405–1424. https://doi.org/10. 1137/1.9781611975482.86
- [3] Deeksha Adil and Sushant Sachdeva. 2020. Faster p-norm minimizing flows, via smoothed q-norm problems. In Proceedings of the 2020 ACM-SIAM Symposium on Discrete Algorithms, SODA 2020, Salt Lake City, UT, USA, January 5-8, 2020, Shuchi Chawla (Ed.). SIAM, 892–910. https://doi.org/10.1137/1.9781611975994.54

- [4] J. Bourgain, J. Lindenstrauss, and V. Milman. 1989. Approximation of zonoids by zonotopes. Acta Math. 162, 1-2 (1989), 73–141. https://doi.org/10.1007/ BF02392835
- [5] Jan van den Brand, Yin Tat Lee, Yang P. Liu, Thatchaphol Saranurak, Aaron Sidford, Zhao Song, and Di Wang. 2021. Minimum cost flows, MDPs, and \$\ell\_1\$-regression in nearly linear time for dense instances. In \$STOC '21: 53rd Annual ACM SIGACT Symposium on Theory of Computing, Virtual Event, Italy, June 21-25, 2021, Samir Khuller and Virginia Vassilevska Williams (Eds.). ACM, 859-869. https://doi.org/10.1145/3406325.3451108
- [6] Jan van den Brand, Yin Tat Lee, Danupon Nanongkai, Richard Peng, Thatchaphol Saranurak, Aaron Sidford, Zhao Song, and Di Wang. 2020. Bipartite Matching in Nearly-linear Time on Moderately Dense Graphs. In 61st IEEE Annual Symposium on Foundations of Computer Science, FOCS 2020, Durham, NC, USA, November 16-19, 2020, Sandy Irani (Ed.). IEEE, 919–930. https://doi.org/10.1109/FOCS46700.2020.00090
- [7] Jan van den Brand, Yin Tat Lee, Aaron Sidford, and Zhao Song. 2020. Solving tall dense linear programs in nearly linear time. In Proceedings of the 52nd Annual ACM SIGACT Symposium on Theory of Computing, STOC 2020, Chicago, IL, USA, June 22-26, 2020, Konstantin Makarychev, Yury Makarychev, Madhur Tulsiani, Gautam Kamath, and Julia Chuzhoy (Eds.). ACM, 775-788. https://doi.org/10. 1145/3357713.3384309
- [8] Sébastien Bubeck, Michael B. Cohen, Yin Tat Lee, and Yuanzhi Li. 2018. An homotopy method for ℓ<sub>p</sub> regression provably beyond self-concordance and in input-sparsity time. In Proceedings of the 50th Annual ACM SIGACT Symposium on Theory of Computing, STOC 2018, Los Angeles, CA, USA, June 25-29, 2018, Ilias Diakonikolas, David Kempe, and Monika Henzinger (Eds.). ACM, 1130–1137. https://doi.org/10.1145/3188745.3188776
- [9] Kenneth L. Clarkson, Ruosong Wang, and David P. Woodruff. 2019. Dimensionality Reduction for Tukey Regression. In Proceedings of the 36th International Conference on Machine Learning, ICML 2019, 9-15 June 2019, Long Beach, California, USA (Proceedings of Machine Learning Research, Vol. 97), Kamalika Chaudhuri and Ruslan Salakhutdinov (Eds.). PMLR, 1262–1271. http://proceedings.mlr.press/v97/clarkson19a.html
- [10] Kenneth L. Clarkson and David P. Woodruff. 2015. Input Sparsity and Hardness for Robust Subspace Approximation. In IEEE 56th Annual Symposium on Foundations of Computer Science, FOCS 2015, Berkeley, CA, USA, 17-20 October, 2015, Venkatesan Guruswami (Ed.). IEEE Computer Society, 310-329. https://doi.org/10.1109/FOCS. 2015.27
- [11] Michael B. Cohen, Yin Tat Lee, and Zhao Song. 2021. Solving Linear Programs in the Current Matrix Multiplication Time. J. ACM 68, 1 (2021), 3:1–3:39. https://doi.org/10.1145/3424305
- [12] Michael B. Cohen and Richard Peng. 2015. Lp Row Sampling by Lewis Weights. In Proceedings of the Forty-Seventh Annual ACM on Symposium on Theory of Computing, STOC 2015, Portland, OR, USA, June 14-17, 2015, Rocco A. Servedio and Ronitt Rubinfeld (Eds.). ACM, 183–192. https://doi.org/10.1145/2746539.2746567
- [13] Mehrdad Ghadiri, Richard Peng, and Santosh S. Vempala. 2021. Faster p-Norm Regression Using Sparsity. abs/2109.11537 (2021). arXiv:2109.11537 https://arxiv. org/abs/2109.11537
- [14] Arun Jambulapati, James R. Lee, Yang P. Liu, and Aaron Sidford. 2023. Sparsifying Sums of Norms. In 64th IEEE Annual Symposium on Foundations of Computer Science, FOCS 2023, Santa Cruz, CA, USA, November 6-9, 2023. IEEE, 1953–1962. https://doi.org/10.1109/FOCS57990.2023.00119
- [15] Arun Jambulapati, Yang P. Liu, and Aaron Sidford. 2022. Improved iteration complexities for overconstrained p-norm regression. In STOC '22: 54th Annual ACM SIGACT Symposium on Theory of Computing, Rome, Italy, June 20 - 24, 2022, Stefano Leonardi and Anupam Gupta (Eds.). ACM, 529-542. https://doi.org/10. 1145/3519935.3519971
- [16] Shunhua Jiang, Zhao Song, Omri Weinstein, and Hengjie Zhang. 2021. A faster algorithm for solving general LPs. In STOC '21: 53rd Annual ACM SIGACT Symposium on Theory of Computing, Virtual Event, Italy, June 21-25, 2021, Samir Khuller and Virginia Vassilevska Williams (Eds.). ACM, 823-832. https://doi.org/10.1145/ 3406325.3451058
- [17] Michel Ledoux and Michel Talagrand. 2011. Probability in Banach spaces. Springer-Verlag, Berlin. xii+480 pages. Isoperimetry and processes, Reprint of the 1991 edition.
- [18] Yin Tat Lee, Zhao Song, and Qiuyi Zhang. 2019. Solving Empirical Risk Minimization in the Current Matrix Multiplication Time. In Conference on Learning Theory, COLT 2019, 25-28 June 2019, Phoenix, AZ, USA (Proceedings of Machine Learning Research, Vol. 99), Alina Beygelzimer and Daniel Hsu (Eds.). PMLR, 2140–2157. http://proceedings.mlr.press/v99/lee19a.html
- [19] D. R. Lewis. 1979. Ellipsoids defined by Banach ideal norms. Mathematika 26, 1 (1979), 18–29. https://doi.org/10.1112/S0025579300009566
- [20] Cameron Musco, Christopher Musco, David P. Woodruff, and Taisuke Yasuda. 2022. Active Linear Regression for \(\ell\_p\) Norms and Beyond. In 63rd IEEE Annual Symposium on Foundations of Computer Science, FOCS 2022, Denver, CO, USA, October 31 - November 3, 2022. IEEE, 744-753. https://doi.org/10.1109/FOCS54457. 2022.00076

- [21] Alain Pajor and Nicole Tomczak-Jaegermann. 1985. Remarques sur les nombres d'entropie d'un opérateur et de son transposé. C. R. Acad. Sci. Paris Sér. I Math. 301, 15 (1985), 743–746.
- [22] Gideon Schechtman and Artem Zvavitch. 2001. Embedding subspaces of  $L_p$  into  $l_p^N$ , 0 . Math. Nachr. 227 (2001), 133–142. https://doi.org/10.1002/1522-2616(200107)227:1%3C133::AID-MANA133%3E3.0.CO;2-8
- [23] Daniel A. Spielman and Nikhil Srivastava. 2011. Graph Sparsification by Effective Resistances. SIAM J. Comput. 40, 6 (2011), 1913–1926.
- [24] Daniel A. Spielman and Shang-Hua Teng. 2014. Nearly Linear Time Algorithms for Preconditioning and Solving Symmetric, Diagonally Dominant Linear Systems. SIAM J. Matrix Anal. Appl. 35, 3 (2014), 835–885. https://doi.org/10.1137/ 090771430
- [25] Michel Talagrand. 1990. Embedding subspaces of  $L_1$  into  $\ell_1^N$ . Proc. Amer. Math. Soc. 108, 2 (1990), 363–369.
- [26] M. Talagrand. 1995. Embedding subspaces of  $L_p$  in  $l_p^N$ . In Geometric aspects of functional analysis (Israel, 1992–1994). Oper. Theory Adv. Appl., Vol. 77.

- Birkhäuser, Basel, 311-325.
- [27] Michel Talagrand. 2014. Upper and lower bounds for stochastic processes. Ergebnisse der Mathematik und ihrer Grenzgebiete. 3. Folge. A Series of Modern Surveys in Mathematics [Results in Mathematics and Related Areas. 3rd Series. A Series of Modern Surveys in Mathematics], Vol. 60. Springer, Heidelberg. xvi+626 pages. https://doi.org/10.1007/978-3-642-54075-2 Modern methods and classical problems.
- [28] David P. Woodruff and Taisuke Yasuda. 2023. Sharper Bounds for \(\ell\_p\) Sensitivity Sampling. In International Conference on Machine Learning, ICML 2023, 23-29 July 2023, Honolulu, Hawaii, USA (Proceedings of Machine Learning Research, Vol. 202), Andreas Krause, Emma Brunskill, Kyunghyun Cho, Barbara Engelhardt, Sivan Sabato, and Jonathan Scarlett (Eds.). PMLR, 37238-37272. https://proceedings.mlr.press/v202/woodruff23a.html

Received 13-NOV-2023; accepted 2024-02-11