

Iterative Refinement for ℓ_p -norm Regression

Deeksha Adil* Rasmus Kyng† Richard Peng‡ Sushant Sachdeva §

January 3, 2019

Abstract

We give improved algorithms for the ℓ_p -regression problem, $\min_{\mathbf{x}} \|\mathbf{x}\|_p$ such that $\mathbf{Ax} = \mathbf{b}$, for all $p \in (1, 2) \cup (2, \infty)$. Our algorithms obtain a high accuracy solution in $\tilde{O}_p(m^{\frac{|p-2|}{2p+|p-2|}}) \leq \tilde{O}_p(m^{1/3})$ iterations, where each iteration requires solving an $m \times m$ linear system, with m being the dimension of the ambient space.

Incorporating a procedure for maintaining an approximate inverse of the linear systems that we need to solve at each iteration, we give algorithms for solving ℓ_p -regression to $1/\text{poly}(n)$ accuracy that runs in time $\tilde{O}_p(m^{\max\{\omega, 7/3\}})$, where ω is the matrix multiplication constant. For the current best value of $\omega > 2.37$, this means that we can solve ℓ_p regression as fast as ℓ_2 regression, for all constant p bounded away from 1.

Our algorithms can be combined with nearly-linear time solvers for linear systems in graph Laplacians to give minimum ℓ_p -norm flow / voltage solutions to $1/\text{poly}(n)$ accuracy on an undirected graph with m edges in $\tilde{O}_p(m^{1+\frac{|p-2|}{2p+|p-2|}}) \leq \tilde{O}_p(m^{4/3})$ time.

For sparse graphs and for matrices with similar dimensions, our iteration counts and running times improve upon the p -norm regression algorithm by [Bubeck-Cohen-Lee-Li STOC'18], as well as general purpose convex optimization algorithms. At the core of our algorithms is an iterative refinement scheme for ℓ_p -norms, using the quadratically-smoothed ℓ_p -norms introduced in the work of Bubeck *et al.* Formally, given an initial solution, we construct a problem that seeks to minimize a quadratically-smoothed ℓ_p norm over a subspace, such that a crude solution to this problem allows us to improve the initial solution by a constant factor, leading to algorithms with fast convergence.

*University of Toronto. deeksha@cs.toronto.edu. Supported by an Ontario Graduate Scholarship, and by a Connaught New Researcher award to Sushant Sachdeva.

†Harvard. rjkyng@gmail.com. Supported by ONR grant N00014-18-1-2562.

‡Georgia Tech. richard.peng@gmail.com. Supported in part by the National Science Foundation under Grant No. 1718533.

§University of Toronto. sachdeva@cs.toronto.edu. Research supported in part by the Natural Sciences and Engineering Research Council of Canada (NSERC), and a Connaught New Researcher award.

Contents

1	Introduction	1
1.1	Contributions	1
1.2	Comparison to Previous Works	3
2	Technical Overview	4
3	Preliminaries	5
4	Main Iterative Algorithm	6
5	Solving the Residual Problem	9
5.1	Equivalent Problems	10
5.2	Oracle	12
5.3	The Algorithm	14
5.4	Proof of Theorem 5.1	20
6	Speedups for General Matrices via. Inverse Maintenance	21
7	Other Regression Formulations	29
7.1	Affine transformations within the norm	29
7.2	$1 < p < 2$	30
8	p-Norm Optimization on Graphs	32
8.1	p -Norm Flows	33
8.2	Lipschitz Learning and Graph Labelling	34
A	Missing Proofs	38
A.1	Proofs from Section 3	38
A.2	Proofs from Section 4	41
A.3	Proofs from Section 5	43
A.4	Proofs from Section 6	48
B	Controlling Φ	49
C	Solving L2 problems	51
D	General ℓ_2 Resistance Monotonicity	52

1 Introduction

Iterative methods that converge rapidly to a solution are of fundamental importance to numerical analysis, optimization, and more recently, graph algorithms. In the study of iterative methods, there are significant discrepancies between iterative methods geared towards linear problems, and ones that can handle more general convex objectives. For systems of linear equations, which corresponds to minimizing ℓ_2 -norm objectives over a subspace, most iterative methods obtain ϵ -approximate solutions in iteration counts that scale as $\log(1/\epsilon)$. More generally, for appropriately defined notions of accuracy, a constant-accuracy linear system solver can be iterated to give a much higher accuracy solver using a few calls to the crude solver. Such phenomena are not limited to linear systems either: an algorithm that produces approximate maximum flows on directed graphs can be iterated on the residual graph to quickly obtain high-accuracy answers.

On the other hand, for the much wider space of non-linear optimization problems arising from optimization and machine learning, it's significantly more expensive to obtain high accuracy solutions. Many widely used methods such as (accelerated) gradient descent, obtain ϵ -approximate answers using iteration counts that scale as $\text{poly}(1/\epsilon)$. Such discrepancies also occur in the overall asymptotic running times. An important and canonical problem in this space is ℓ_p -norm regression:

$$\min_{\mathbf{x} \in \mathbb{R}^m: \mathbf{A}\mathbf{x}=\mathbf{b}} \|\mathbf{x}\|_p^p, \quad (*)$$

for some $\mathbf{A} \in \mathbb{R}^{n \times m} (m \geq n)$, and $\mathbf{b} \in \mathbb{R}^n$. For $p = 2$, this corresponds exactly to solving a linear system, and hence is solvable by a matrix inversion in $O(m^\omega)$ time ¹ For $p = 1$ and $p = \infty$, this problem is inter-reducible to linear programming [5, 42, 43].

Interior point methods also allow us to solve ℓ_p -norm regression problems in $\sqrt{\text{rank}}$ iterations [26, 31], where each iteration requires solving an $m \times m$ linear system for any $p \in [1, \infty]$. Bubeck *et al* [5] show that this iteration count is tight for the interior point method framework, and instead propose a different method which requires only $\tilde{O}_p(m^{|\frac{1}{2}-\frac{1}{p}|})$ ² iterations for $p \neq 1, \infty$, which for large constant p still tends to about $m^{1/2}$. On the other hand, ϵ -approximate solutions can be computed in about $m^{1/3} \text{poly}(1/\epsilon)$ iterations [6] ³.

Furthermore, this discrepancy also carries over to the graph theoretic case. If the matrix \mathbf{A} is the vertex-edge incidence matrix of a graph, then this problem captures graph problems such as p -norm Lipschitz learning and finding ℓ_p -norm minimizing flows meeting demands given by \mathbf{b} . Here low accuracy approximate solutions can be obtained in nearly-linear time when $p = \infty$ [33, 36], and almost-linear time for all other values of p [37, 39]. However, the current best high accuracy solutions take at least $m^{\min\{10/7, 1+|1/2-1/p|\}}$ time [5, 28].

1.1 Contributions

Iterative Refinement for ℓ_p -norms. In this paper, we propose a new iterative method for ℓ_p -norm regression problems (*) that achieves geometric convergence to the optimal solution. Our method only requires solving $O_p(\log 1/\epsilon)$ *residual problems* to find an ϵ -approximate solution, or $O_p(\kappa \log 1/\epsilon)$ residual problems, each solved to a κ -approximation factor. Such an iterative method

¹ ω is the matrix multiplication exponent. Currently we know $\omega \leq 2.3728639$. [22, 46].

² $O_p(\cdot)$ notation hides constant factors that depend on p , and its dual norm $\frac{p}{p-1} \cdot \tilde{O}_p(\cdot)$ notation also hides $\text{poly}(\log \frac{mn}{\epsilon})$ factors in addition.

³this result only addressed the $p = \infty$ case, but its techniques generalize to all other p

was previously known only for $p = 2$ and ∞ . Curiously, our residual problems look very similar to the original problem (*), with the ℓ_p norms replaced by their quadratically-smoothed versions introduced by Bubeck *et al* [5]. This result, Theorem 4.1, can be stated informally as:

Theorem 1.1. *There exists a class of residual problems for p -norm regression (which we will define in Definition 4.3) such that any p -norm regression problem can be solved to ϵ -relative accuracy by solving to relative error κ a sequence of $O_p(\kappa \log(\frac{m}{\epsilon}))$ residual problems.*

Improved Iteration Count for ℓ_p -Regression. We then give an algorithm for quickly solving the residual problem motivated by the approximate maximum flow by electrical flows algorithm by Christiano *et al.* [7] and its generalizations to regression problems [6]. This is given as Theorem 5.1, and can be stated informally as:

Theorem 1.2. *For any $p > 2$, an instance of a residual problem for p -norm regression as defined in Definition 4.3 can be solved in $\tilde{O}_p(m^{\frac{p-2}{3p-2}})$ iterations, each of which consist of solving a system of linear equations plus updates that take linear time.*

This improves on the work of Bubeck *et al* [5] for all $p > 2$, with the number of iterations equaling $\tilde{O}_p(1)$ for $p = 2$ (essentially the same as Bubeck *et al*) and tending to $\tilde{O}(m^{1/3})$ as p goes to ∞ (compared to $\tilde{O}(m^{1/2})$ for Bubeck *et al*). However, our results don't give anything for $p = \infty$ due to the dependency in p in the $\tilde{O}_p(\cdot)$ term. It's worth noting that even in the constant error regime, this improves by a factor of about $\min\{m^{\frac{(p-2)^2}{2p(3p-2)}}, m^{\frac{4}{3(3p-2)}}\}$ over the current state of the art, which for small p is due to Bubeck *et al.* [5], and for large p is based on unpublished modifications to Christiano *et al.* [7, 29].

A Duality Based Approach to ℓ_p -Regression. For the remaining case of $1 < p < 2$, we instead solve the dual problem, which is a $\frac{p}{p-1}$ -norm regression problem, and utilize its solution to solve our original ℓ_p -Regression. This leads to iteration counts of the form $\tilde{O}_p(m^{\frac{2-p}{p+2}} \log(1/\epsilon))$ for solving such problems. Note that this result also does not give anything when $p = 1$, as the constants related to its dual norm, $\frac{p}{1-p}$ become prohibitive. For all $p \in (1, \infty)$, our iteration count achieves the following exponent on m ,

$$\frac{\left|\frac{1}{2} - \frac{1}{p}\right|}{1 + \left|\frac{1}{2} - \frac{1}{p}\right|},$$

while the exponent from the previous result [5] is $\left|\frac{1}{2} - \frac{1}{p}\right|$: our algorithm has better dependence on m on all constant p (albeit with larger constants depending on p).

Faster ℓ_p -Regression and p -Norm Flows. Our improved iteration counts can be readily combined with methods for speeding up algorithms that utilize linear system solvers, including fast Laplacian solvers for graph problems [41, 45] and inverse maintenance [25, 44]. This results in:

1. Algorithms for finding p -norm flow on undirected graphs to accuracy ϵ with running time $\tilde{O}_p\left(m^{1 + \frac{|p-2|}{2p+|p-2|}} \log(1/\epsilon)\right)$ for $p \in (1, \infty)$ via direct invocations of fast Laplacian solvers [40].

2. An $\tilde{O}_p(m^{\max\{\omega, 7/3\}})$ time algorithm for solving ℓ_p regression problems for all $p \in (1, \infty)$, which we formalize in Theorem 6.1.

The latter bound for general matrices brings us to the somewhat surprising conclusion that for the current value of $\omega > 7/3$, p -norm regression problems (with constant p that's also constant-bounded away from 1) on square matrices can be solved as fast as solving the underlying linear systems, or equivalently, ℓ_2 regression problems.

The latter of these maintains an approximate inverse to the linear systems we need to solve in each step of the iterative method as pioneered by Vaidya [44]. However, our modification interacts directly with the potential functions we use to control iteration counts in the inner loop of our iterative method. An approach based on similar principles was also given by Cohen et al., with announcement [24] and manuscript availability [8] taking place before and after the initial submission of this paper respectively.

When combined with fast solvers for linear systems in graph Laplacians, our results give the first evidence that wide classes of optimization problems can be solved in time $m^{4/3}$ or faster. While such a bound (via. fast Laplacian solvers) is by now well-known in the approximate setting [7], the $m^{10/7}$ iteration bounds due to Madry [28, 30] represent the only results to date in this direction for high accuracy answers on sparse graphs.

Finally we remark that for the case of matrices with far more rows than columns, as well as dense graphs, our running times are worse than the input-sparsity time bounds from Bubeck et al. [5]. However, we believe that the (input-sparsity time) p -norm sampling algorithms from Cohen and Peng [11] should be applicable to our residual problems, yielding input-sparsity time variants of our algorithms. However, an analysis of their routines for our residual problems containing mixed ℓ_2 and ℓ_p norm functions is outside the scope of this paper.

Generalizations and Extensions. While we focus on Problem (*), under mild assumptions about polynomially bounded objectives, we can solve the following more general problem:

$$\begin{aligned} \min_x \quad & \|C\mathbf{x} - \mathbf{d}\|_p \\ & A\mathbf{x} = \mathbf{b} \end{aligned}$$

The reduction is discussed in Section 7. The combination of an affine constraint on \mathbf{x} with an affine transformation in the p -norm objective means we can solve most variants of p -norm optimization problems.

Similar ideas can be used to solve p -norm Lipschitz learning problems [21] on graphs quickly.

1.2 Comparison to Previous Works

Numerical Methods and Preconditioning. Iterative methods and preconditioning are the most fundamental tools in numerical algorithms [4, 35]. As studies of such methods often focus on linear problems, many existing analyses of iterative methods are restricted to linear systems. Generalizing such methods, as well as numerical methods, to broader settings is a major topic of study [15, 16, 19, 32].

The study of more efficient algorithms for combinatorial flow problems has benefited enormously from ideas from linear and non-linear preconditioning. Recent advances in approximate maximum flow and transshipment algorithms [13, 17, 23, 33, 34, 36, 37, 38] build upon such ideas. However, these methods rely on the preconditioner being a linear operator, and give $poly(1/\epsilon)$ dependence.

Optimization Algorithms. Our techniques for solving the residual problems are directly motivated by approximating maximum flow using electrical flows [7]. While this algorithm has been extended to multicommodity flows and regression problems [6, 18], all these results have $poly(1/\epsilon)$ dependencies.

Several recent results for obtaining $\log(1/\epsilon)$ dependencies are all motivated by convex optimization techniques. In particular, the state of the art running times are by interior point methods. These include directly modifying the interior point method (IPM) [21, 25, 26], combining techniques from the electrical flow algorithms with IPM update steps [9, 20, 21, 28, 30], and increasing the ‘confidence interval’, and in turn step lengths, of the IPM update steps [2, 5, 10]. Our result based on creating intermediate problems has the most in common with the last of these. However, our method differs in that our guarantees for this intermediate problem holds over the entire space.

Inverse Maintenance and Numerical Linear Algebra Algorithms. Our final running time of $\tilde{O}_p(m^{\max\{\omega, 7/3\}})$ for ℓ_p -regression incorporates inverse maintenance. This is a method introduced by Vaidya [44] for speeding up optimization algorithms for solving minimum cost and multicommodity flows. It takes advantage of the controllable rate at which optimization algorithms modify the solution variables to reuse inverses of matrices constructed from such variables.

Previous studies of inverse maintenance [25, 26, 44] have been geared towards the interior point method. Here the norm per update step can be controlled, and we believe this also holds for their applications in faster cutting plane methods [27]. While such methods also give gains in the case of our algorithm, for the final bound of about m^ω , we instead bound the progress of the steps against a global potential function motivated by the electrical flow max-flow algorithm [7].

Our algorithm on the other hand does not take into account sparsity of the input matrix, or possibly uneven dimensions (e.g. $m \approx n^2$). In these settings, the methods based on accelerated stochastic gradient descent from [5] obtain better performances. On the other hand, we believe our methods have the potential of extending to such settings by combining the intermediate problems with ℓ_p row sampling [11].

2 Technical Overview

Iterative Refinement for ℓ_p -norms. To design their algorithm for ℓ_p -norm regression, Bubeck *et al* [5] construct a function $\gamma_p(t, x)$, which is C_1 ,⁴ quadratic in the range $|x| \leq t$, and behaves as $|x|^p$ asymptotically (see Def. 3.1). Our key lemma states one can locally approximate $\|x + \Delta\|_p^p$ as a linear function plus a $\gamma_p(|x|, \Delta)$ “error” term⁵ (Lemma 4.5):

$$|x + \Delta|^p = |x|^p + \Delta \frac{d}{dx} |x|^p + O_p(1) \gamma_p(|x|, \Delta).$$

Surprisingly, this approximation only has an $O_p(1)$ “condition number”. Proceeding just as for gradient descent, or Newton’s method, means that if at each step we solve the following local approximation problem to a factor κ ,

$$\max_{\Delta} \mathbf{g}^\top \Delta - \gamma_p(\mathbf{x}, \Delta),$$

⁴A function is said to be C_1 if it’s continuous, differentiable, and has a continuous derivative

⁵It is useful to compare the γ_p term to the second-order Hessian term in Taylor expansion

where \mathbf{g} is the gradient of our loss function, we can converge to an ε -approximate solution in roughly $\kappa \log 1/\varepsilon$ iterations (Theorem 4.1).

Improved Algorithms for ℓ_p -regression for $p \geq 2$. The key advantage afforded by our iterative algorithm is that we now only need to design an algorithm for the residual problem that achieves a crude approximation factor (we achieve $O_p(1)$). As a first step, by a binary search and some rescaling, we show (Lemma 5.5) that it suffices to achieve a constant factor approximation to $O_p(\log m/\varepsilon)$ problems of the following form,

$$\min_{\mathbf{A}\mathbf{x}=0, \mathbf{g}^\top \mathbf{x}=c} \gamma_p(\mathbf{t}, \mathbf{x}).$$

The technical heart of our proof is to give an algorithm (Gamma-Solver, Algorithm 4) inspired by the multiplicative weight update (MWU) method (see [3] for a survey), combined with the width-reduction inspired by the faster flow algorithm of Christiano *et al.* [7], and its matrix version by Chin *et al.*[6]. At each iteration, we solve a weighted ℓ_2 minimization problem to find the next update step. If this update step has small ℓ_p norm, we add this to our current solution, and update the weights. Otherwise, we identify a set of coordinates that have small current weights, and yet are contributing most of the ℓ_p norm, and we penalize them by increasing their weights (and do not add our update step to the current solution). Setting the parameters carefully, we show that after $\tilde{O}_p(m^{\frac{p-2}{3p-2}})$ iterations, the average of the update steps achieves an $O_p(1)$ -approximation to our modified residual problem (Theorem 5.8). Combining this with our iterative refinement algorithm, we obtain our algorithms for ℓ_p -norm regression that require only $\tilde{O}_p(m^{\frac{p-2}{3p-2}})$ iterations (or linear system solves).

Maintaining Inverses for Improved Algorithm. Our inverse maintenance procedure utilizes the same combination of low-rank updates and matrix multiplications as in previous results [25, 26, 44]. However, the rate of convergence of our algorithm, and in turn the rate at which we adjust the weights from the MWU procedure, are governed by growths in the ℓ_2 minimization problem. This leads to the difficulty of uneven progress across the iterations.

We solve this issue by a simple yet subtle scheme motivated by lazy updates in data structures [1, 14]. We bucket changes to the values of entries based on their magnitudes, and update entries that received too many updates of a certain magnitude separately. This differs with previous methods that update weights exceeding approximation thresholds as they happen, and enables a closer interaction with the overall potential function based convergence analysis.

3 Preliminaries

We use the following family of functions, $\gamma_p(t, x)$ defined in [5].

Definition 3.1 (γ_p function). For $t \geq 0$ and $p \geq 1$, define

$$\gamma_p(t, x) = \begin{cases} \frac{p}{2}t^{p-2}x^2 & \text{if } |x| \leq t, \\ |x|^p + (\frac{p}{2} - 1)t^p & \text{otherwise.} \end{cases}$$

These functions can be thought of a quadratic approximation of $|x|^p$ in a small range around zero. The following properties follow directly from the definition.

1. $\gamma_p(0, x) = |x|^p$.
2. $\gamma_p(t, x)$ is quadratic in the range $-t \leq x \leq t$.
3. γ_p is C^1 in both x, t .

We show several other important properties of γ_p in the following lemmas. Their proofs are straightforward and deferred to Appendix A.1

Lemma 3.2. *Function γ_p is as defined above.*

1. For any $p \geq 2$, $t \geq 0$, and $x \in \mathbb{R}$, we have $\gamma_p(t, x) \geq |x|^p$, and $\gamma_p(t, x) \geq \frac{p}{2}t^{p-2}x^2$.
2. It is homogeneous under rescaling of both t and x , i.e., for any $t, \lambda \geq 0, p \geq 1$, and any x we have $\gamma_p(\lambda t, \lambda x) = \lambda^p \gamma_p(t, x)$.
3. For any $t > 0, p \geq 1$ and any x , we have $\gamma_p'(t, x) = p \max\{t, |x|\}^{p-2}x$.

The next lemma shows a bound on the value of γ_p when x is scaled up or down.

Lemma 3.3. *For any $p > 1, \Delta \in \mathbb{R}$ and $\lambda \geq 0$, we have,*

$$\min\{2, p\} \leq x \frac{\gamma_p'(t, x)}{\gamma_p(t, x)} \leq \max\{2, p\}.$$

This implies,

$$\min\{\lambda^2, \lambda^p\} \gamma_p(t, \Delta) \leq \gamma_p(t, \lambda \Delta) \leq \max\{\lambda^2, \lambda^p\} \gamma_p(t, \Delta).$$

The following lemma allows us to bound the second order change in $\gamma_p(\mathbf{x})$ as \mathbf{x} changes to $\mathbf{x} + \Delta$.

Lemma 3.4. *For any $p \geq 2, t \geq 0$ and any x, Δ , we have*

$$\gamma_p(t, x + \Delta) \leq \gamma_p(t, x) + \left| \gamma_p'(t, x) \Delta \right| + p^2 2^{p-3} \max\{t, |x|, |\Delta|\}^{p-2} \Delta^2.$$

Notation. For a vector \mathbf{x} , let $|\mathbf{x}|$ denote the vector with its i^{th} coordinate as $|x_i|$. For any two vectors \mathbf{t} and \mathbf{x} , $\gamma_p(\mathbf{t}, \mathbf{x})$ denotes the sum $\sum_i \gamma_p(t_i, x_i)$.

4 Main Iterative Algorithm

In this section we analyze procedure P-NORM, i.e., Algorithm 1. Our main result for this section is,

Theorem 4.1 (ℓ_p -norm Iterative Refinement). *For any $p \in (1, \infty)$, and $\kappa \geq 1$. Given an initial feasible solution $\mathbf{x}^{(0)}$ (Definition 4.7) to our optimization problem (Equation (*)), Algorithm 1 finds an ϵ -approximate solution to (*) in $O_p\left(\kappa \log\left(\frac{m}{\epsilon}\right)\right)$ calls to a κ -approximate solver for the residual problem (Equation (1)).*

The theorem says that it is sufficient to solve an instance of the residual problem (1) crudely, and only a logarithmic number of times. Before we prove the theorem, we define the terms used in the statement and prove some results that would be needed for the proof. We begin by defining an ϵ -approximate solution to our main optimization problem.

Algorithm 1 Meta-Iterative Algorithm

```

1: procedure P-NORM( $\mathbf{A}, \mathbf{b}, \varepsilon$ )
2:    $\mathbf{x}^{(0)} \leftarrow \min_{\mathbf{A}\mathbf{x}=\mathbf{b}} \|\mathbf{x}\|_2^2$ .
3:    $T \leftarrow O_p(\kappa \log(\frac{m}{\varepsilon}))$ 
4:    $\lambda \leftarrow \left(\frac{p-1}{p^{4p}}\right)^{\frac{1}{\min(1, p-1)}}$ 
5:   for  $t = 0$  to  $T$  do
6:      $\Delta \leftarrow \kappa\text{-APPROX}(\mathbf{A}, \mathbf{b}, \mathbf{x}^{(t)}, \lambda, \mathbf{x}^{(0)})$   $\triangleright \kappa$ -approximate solution to (1)
7:      $\mathbf{x}^{(t+1)} \leftarrow \mathbf{x}^{(t)} - \lambda\Delta$ 
8:     if  $\|\mathbf{x}^{(t+1)}\|_p^p \geq \|\mathbf{x}^{(t)}\|_p^p$  then
9:       return  $\mathbf{x}^{(t)}$ 
10:  return  $\mathbf{x}^{(T)}$ 

```

Definition 4.2 (ε -approximate solution). *We say our solution \mathbf{x} is an ε -approximate solution to (*) if $\mathbf{A}\mathbf{x} = \mathbf{b}$ and*

$$\|\mathbf{x}\|_p^p \leq (1 + \varepsilon)\|\mathbf{x}^*\|_p^p,$$

where \mathbf{x}^* is the OPT of (*).

We next define what we use as our residual problem and what we mean by a κ -approximate solution.

Definition 4.3 (Residual Problem). *For any given \mathbf{x} and $p > 1$, let*

$$\alpha(\Delta) \stackrel{\text{def}}{=} \langle \mathbf{g}, \Delta \rangle - \frac{p-1}{p^{2p}} \gamma_p(|\mathbf{x}|, \Delta),$$

where \mathbf{g} is the gradient, $\mathbf{g} = p|\mathbf{x}|^{p-2}\mathbf{x}$. We call the following problem to be the residual problem of (*) at \mathbf{x} .

$$\max_{\mathbf{A}\Delta=0} \alpha(\Delta). \tag{1}$$

Definition 4.4 (κ -approximate solution). *Let $\kappa \geq 1$. A κ -approximate solution for the residual problem is $\tilde{\Delta}$ such that $\mathbf{A}\tilde{\Delta} = 0$ and, $\alpha(\tilde{\Delta}) \geq \frac{1}{\kappa}\alpha(\Delta^*)$. Here $\Delta^* = \max_{\mathbf{A}\Delta=0} \alpha(\Delta)$.*

In order to see why we choose this problem as our residual problem we show that the objective of the residual problem bounds the change in p -norm of a vector \mathbf{x} when perturbed by Δ (Lemma 4.6).

Lemma 4.5. *Let $p \in (1, \infty)$. Then for any \mathbf{x} and any Δ ,*

$$|\mathbf{x}|^p + g\Delta + \frac{p-1}{p^{2p}} \gamma_p(|\mathbf{x}|, \Delta) \leq |\mathbf{x} + \Delta|^p \leq |\mathbf{x}|^p + g\Delta + 2^p \gamma_p(|\mathbf{x}|, \Delta),$$

where $g = p|\mathbf{x}|^{p-2}\mathbf{x}$ is the derivative of the function $|\mathbf{x}|^p$.

The proof can be found in Appendix A.2.

Lemma 4.6. *Let $p \in (1, \infty)$ and λ be such that $\lambda^{\min\{1, p-1\}} \leq \frac{p-1}{p^{4p}}$. Then for any Δ we have,*

$$\|\mathbf{x}\|_p^p - \alpha(\lambda\Delta) \leq \|\mathbf{x} - \lambda\Delta\|_p^p \leq \|\mathbf{x}\|_p^p - \lambda\alpha(\Delta).$$

Proof. Applying lemma 4.5 to all the coordinates, we obtain,

$$\|\mathbf{x}\|_p^p - \langle g, \Delta \rangle + \frac{p-1}{p2^p} \gamma_p(|\mathbf{x}|, \Delta) \leq \|\mathbf{x} - \Delta\|_p^p \leq \|\mathbf{x}\|_p^p - \langle g, \Delta \rangle + 2^p \gamma_p(|\mathbf{x}|, \Delta).$$

Using definition 4.3, equation (4) directly implies, $\|\mathbf{x} - \Delta\|_p^p \geq \|\mathbf{x}\|_p^p - \alpha(\Delta)$ for all Δ . Now to prove the other side, note that for any $\lambda \in [0, 1]$, and any Δ , we have from Lemma 4.5 and Lemma 3.2

$$\|\mathbf{x} - \lambda\Delta\|_p^p \leq \|\mathbf{x}\|_p^p - \langle g, \lambda\Delta \rangle + 2^p \gamma_p(|\mathbf{x}|, \lambda\Delta) \leq \|\mathbf{x}\|_p^p - \lambda \left(\langle g, \Delta \rangle - \lambda^{\min\{1, p-1\}} 2^p \gamma_p(|\mathbf{x}|, \Delta) \right).$$

Picking λ such that $\lambda^{\min\{1, p-1\}} \leq \frac{p-1}{p4^p}$, we obtain that for any Δ ,

$$\|\mathbf{x} - \lambda\Delta\|_p^p \leq \|\mathbf{x}\|_p^p - \lambda \left(\langle g, \Delta \rangle - \frac{p-1}{p2^p} \gamma_p(|\mathbf{x}|, \Delta) \right) = \|\mathbf{x}\|_p^p - \lambda\alpha(\Delta),$$

thus concluding the proof of the lemma. \square

For any iterative algorithm we need a starting feasible solution. We could potentially start with any feasible solution but we define the following starting solution which we claim is a *good* starting point. Lemma 4.8 shows us that our chosen starting point is only polynomially away from the optimum solution, and is thus a *good* choice. The proof of the lemma can be found in Appendix A.2.

Definition 4.7 (Initial Solution). *We define $\mathbf{x}^{(0)}$ to be our initial feasible solution to be $\mathbf{x}^{(0)} = \min_{A\mathbf{x}=b} \|\mathbf{x}\|_2^2$.*

Lemma 4.8. *For $\mathbf{x}^{(0)}$ as defined in Definition 4.7, $\|\mathbf{x}^{(0)}\|_p^p \leq m^{(p-2)/2} \text{OPT}$.*

We are now ready to prove Theorem 4.1.

Proof. Let $\tilde{\Delta}$ denote the solution returned by the κ -approximate solver. We know that $\alpha(\tilde{\Delta}) \geq \frac{1}{\kappa} \cdot \alpha(\Delta^*)$. We have,

$$\alpha(\tilde{\Delta}) \geq \frac{1}{\kappa} \cdot \alpha(\Delta^*) \geq \frac{1}{\kappa} \alpha(\mathbf{x} - \mathbf{x}^*) \geq \frac{1}{\kappa} \left(\|\mathbf{x}\|_p^p - \|\mathbf{x}^*\|_p^p \right) = \frac{1}{\kappa} \left(\|\mathbf{x}\|_p^p - \text{OPT} \right).$$

From Lemma 4.6, for $\lambda = \left(\frac{p-1}{p4^p} \right)^{\frac{1}{\min\{1, p-1\}}} = \Omega_p(1)$, we get,

$$\left\| \mathbf{x} - \lambda\tilde{\Delta} \right\|_p^p \leq \|\mathbf{x}\|_p^p - \lambda\alpha(\tilde{\Delta}).$$

Combining the above two equations and subtracting OPT from both sides gives us

$$\begin{aligned} \left\| \mathbf{x} - \lambda\tilde{\Delta} \right\|_p^p - \text{OPT} &\leq -\lambda\alpha(\tilde{\Delta}) + \|\mathbf{x}\|_p^p - \text{OPT} \\ &\leq -\frac{\lambda}{\kappa} \left(\|\mathbf{x}\|_p^p - \text{OPT} \right) + \left(\|\mathbf{x}\|_p^p - \text{OPT} \right) \\ &\leq \left(1 - \frac{\lambda}{\kappa} \right) \left(\|\mathbf{x}\|_p^p - \text{OPT} \right). \end{aligned}$$

Using lemma 4.8 we get,

$$\mathbf{x}^{(t)} - \text{OPT} \leq \left(1 - \frac{\lambda}{\kappa}\right)^t \left(\mathbf{x}^{(0)} - \text{OPT}\right) \leq \left(1 - \frac{\lambda}{\kappa}\right)^t \left(m^{\frac{p-2}{2}} - 1\right) \text{OPT}.$$

Setting $t = O_p\left(\kappa \log\left(\frac{m}{\varepsilon}\right)\right)$ gives us an ε -approximate solution. \square

This concludes the discussion on the analysis of Algorithm 1. In the following sections we move on to analyzing how to solve the residual problem approximately.

5 Solving the Residual Problem

In this section, we give an algorithm that solves the residual problem to a constant approximation. Combined with the iterative refinement scheme from Theorem 4.1, we obtain the following result.

Theorem 5.1. *For $p \geq 2$, we can find an ε -approximate solution to (*) in time*

$$\tilde{O}_p\left((m+n)^{\omega + \frac{p-2}{3p-2}} \log^2 1/\varepsilon\right).$$

Here ω is the matrix multiplication constant.

Recall that the residual problem

$$\max_{\mathbf{A}\Delta=0} \mathbf{g}^\top \Delta - \frac{p-1}{p2^p} \gamma_p(\mathbf{t}, \Delta),$$

has a linear term followed by the γ_p function. Instead of directly optimizing this function, we guess an approximate value of the linear term, and for each such guess, we minimize the γ_p function under this additional constraint. We can scale the problem so that the optimum is at most 1. Finally, we can perturb \mathbf{t} so that each t_i lies in a polynomially bounded range without adding significant error. Our final problem looks as follows,

$$\begin{aligned} \min_{\Delta} \quad & \gamma_p(\mathbf{t}, \Delta) \\ & \mathbf{A}\Delta = 0, \\ & \mathbf{g}^\top \Delta = c, \end{aligned} \tag{2}$$

with $m^{-1/p} \leq t_i \leq 1, \forall i$.

To summarise, κ -APPROX (Algorithm 2) formalizes this process and shows that we only need to solve a logarithmic number of instances of the above program, (2) and solving each to a κ -approximation gives a $\Omega_p\left(\kappa^{1/(\min\{2,p\}-1)}\right)$ -approximate solution to (1). GAMMA-SOLVER (Algorithm 4) solves problem (2) to an $O_p(1)$ approximation. Therefore, using GAMMA-SOLVER as a subroutine for κ -APPROX we get an $O_p(1)$ approximate solution to (1). Section 5.1 gives an analysis for κ -APPROX. In Section 5.2, we give an oracle that is used in GAMMA-SOLVER. We give an analysis of GAMMA-SOLVER in Section 5.3. Finally in Section 5.4, we give a proof for Theorem 5.1.

5.1 Equivalent Problems

In this section we prove the following theorem.

Theorem 5.2. *Procedure κ -APPROX (Algorithm 2) returns an $\Omega_p \left(\kappa^{1/(\min\{2,p\}-1)} \right)$ -approximate solution to the residual problem given by (1), by solving $O_p \left(\log \left(\frac{m}{\varepsilon} \right) \right)$ instances of program (2) to a κ -approximation.*

The following lemmas will lead to the proof of the above theorem. The first lemma gives an upper and lower bound on the objective of (1).

Lemma 5.3. *Let $p \in (1, \infty)$ and assume that our current solution \mathbf{x} is not an ε -approximate solution. Let λ be such that $\lambda^{\min\{1,p-1\}} = \frac{p-1}{p^4}$. For some*

$$i \in \left[\log \left(\frac{\varepsilon \|\mathbf{x}^{(0)}\|_p^p}{m^{p-2/2}} \right), \log \left(\frac{\|\mathbf{x}^{(0)}\|_p^p}{\lambda} \right) \right],$$

$\alpha(\Delta^*) \in [2^{i-1}, 2^i]$ where Δ^* is the solution of (1).

We defer the proof to Appendix A.3. Lemma 5.3 suggests that we can divide the range of the objective of our residual problem, α into a logarithmic number of bins and solve a *decision problem* that asks if the optimum belongs to the bin. The lemma guarantees that at least one of the decision problems will be feasible. The following lemma defines the required *decision problems* and shows that solving these to a constant approximation is sufficient to get a constant approximate solution to (1).

Lemma 5.4. *Let $p \in (1, \infty)$. Suppose $\alpha(\Delta^*) \in [2^{i-1}, 2^i]$ for some i where Δ^* is the solution of (1). The following program is feasible:*

$$\begin{aligned} \gamma_p(\mathbf{t}, \Delta) &\leq \frac{p}{p-1} 2^{i+p}, \\ \mathbf{g}^T \Delta &= 2^{i-1}, \\ \mathbf{A} \Delta &= 0. \end{aligned} \tag{3}$$

If $\Delta(i)$ is a β -approximate solution to program (3) for this choice of i , then, we can pick $\mu \leq 1$ such that the vector $\mu \Delta(i)$ is an $\Omega_p \left(\beta^{\frac{1}{\min\{p,2\}-1}} \right)$ -approximate solution to (1).

The proof can be found in Appendix A.3. We now scale down the objective of (3) so that it is at most 1. The next lemma shows what the scaled down problem looks like and how an approximate solution to the scaled down problem gives an approximate solution to (3). Again the proof of the lemma can be found in Appendix A.3.

Lemma 5.5. *Let $p \in (1, \infty)$. Let i be such that (3) is feasible. Let*

$$\hat{\mathbf{t}}_j = \begin{cases} m^{-1/p} & \left(\frac{p-1}{p} \right)^{1/p} 2^{-i/p-1} \mathbf{t}_j \leq m^{-1/p}, \\ 1 & \left(\frac{p-1}{p} \right)^{1/p} 2^{-i/p-1} \mathbf{t}_j \geq 1, \\ \left(\frac{p-1}{p} \right)^{1/p} 2^{-i/p-1} \mathbf{t}_j & \text{otherwise.} \end{cases}$$

Note that $m^{-1/p} \leq \hat{\mathbf{t}}_j \leq 1$. Then program (2) with $\mathbf{t} = \hat{\mathbf{t}}$, and

$$\mathbf{c} = \left(\frac{2}{p}\right)^{1/2} \left(\frac{p-1}{p}\right)^{1/p} 2^{i(1-\frac{1}{p})-2},$$

has $\text{OPT} \leq 1$. Let Δ^* be a κ -approximate solution to (2). Then, $\Delta = \left(\frac{p}{2}\right)^{1/2} \left(\frac{p}{p-1}\right)^{1/p} 2^{1+i/p} \Delta^*$ is a $\Omega_p(\kappa)$ -approximate solution to (3).

Algorithm 2 Approximate Solver

```

1: procedure  $\kappa$ -APPROX( $\mathbf{A}, \mathbf{b}, \mathbf{x}, \lambda, \mathbf{x}^{(0)}$ )
2:    $\hat{\mathbf{A}} \leftarrow [\mathbf{A}^\top, \mathbf{g}]^\top$ 
3:
4:   for  $i \in \left[ \log \left( \frac{\varepsilon \|\mathbf{x}^{(0)}\|_p^p}{m^{(p-2)/2}} \right), \log \left( \frac{\|\mathbf{x}^{(0)}\|_p^p}{\lambda} \right) \right]$  do
5:      $c(i) \leftarrow \left(\frac{2}{p}\right)^{1/2} \left(\frac{p-1}{p}\right)^{1/p} 2^{i(1-\frac{1}{p})-2}$ 
6:     For every  $e$ ,
7:        $\mathbf{t}_e \leftarrow \max\{m^{-\frac{1}{p}}, (\frac{p-1}{p})^{1/p} 2^{-i/p-1} |\mathbf{x}_e^{(t)}|\}$ 
8:       For every  $e$ ,  $\mathbf{t}_e \leftarrow \min\{1, \mathbf{t}_e\}$ 
9:        $\mathbf{c} \leftarrow [\mathbf{0}^\top, c(i)]^\top$ 
10:       $\Delta^* \leftarrow \text{GAMMA-SOLVER}(\hat{\mathbf{A}}, \mathbf{c}, \mathbf{t})$  ▷  $\kappa$ -approximation to (2)
11:       $\Delta^{(i)} \leftarrow \left(\frac{p}{2}\right)^{1/2} \left(\frac{p}{p-1}\right)^{1/p} 2^{1+i/p} \Delta^*$ 
12:       $\beta \leftarrow \left(\frac{p}{2}\right)^{p/2} \frac{p2^{p+1}}{p-1} 2^i \kappa$ 
13:       $\mu^{(i)} \leftarrow \begin{cases} \left(\frac{1}{2\beta p}\right)^{1/(p-1)} & \text{if } p \leq 2 \\ \frac{1}{4\beta} & \text{otherwise.} \end{cases}$ 
14:       $\Delta^{(i)} \leftarrow \mu^{(i)} \Delta^{(i)}$ 
15:   return  $\lambda \cdot \arg \min_{\Delta^{(i)}} \|\mathbf{x}^{(t)} - \lambda \Delta^{(i)}\|_p^p$ 

```

We now prove Theorem 5.2.

Proof. Lemma 5.3 suggests that there exists an index

$$j \in \left[\log \left(\frac{\varepsilon \|\mathbf{x}^{(0)}\|_p^p}{m^{(p-2)/2}} \right), \log \left(\frac{\|\mathbf{x}^{(0)}\|_p^p}{\lambda} \right) \right],$$

such that $\text{OPT} = \alpha(\Delta^*) \in [2^{j-1}, 2^j]$. Lemma 5.4 implies that (3) is feasible for index j . Suppose $\Delta^{(j)}$ is a κ -approximate solution to the scaled down problem (2) for index j . Lemma 5.5 implies that $\tilde{\Delta}^{(j)} = \left(\frac{p}{2}\right)^{1/2} \left(\frac{p}{p-1}\right)^{1/p} 2^{1+i/p} \Delta^{(j)}$ is an $\Omega_p(\kappa)$ approximate solution to (3) for index j . Lemma 5.4 now implies that $\tilde{\Delta}^{(j)} = \mu \Delta^{(j)}$ is a $\Omega_p \left(\kappa^{\frac{1}{\min\{p, 2\}-1}} \right)$ -approximation to our residual problem

(1). Now, the algorithm solves the scaled down problem for every i and returns the $\tilde{\Delta}^{(i)}$ that when added to our current solution gives the minimum ℓ_p -norm. It either chooses $\tilde{\Delta}^{(j)}$ or some other solution $\tilde{\Delta}^{(i)}$. In case it returns $\tilde{\Delta}^{(i)}$,

$$\begin{aligned} \|\mathbf{x}\|_p^p - \|\mathbf{x} - \lambda\tilde{\Delta}^{(i)}\|_p^p &\geq \|\mathbf{x}\|_p^p - \|\mathbf{x} - \lambda\tilde{\Delta}^{(j)}\|_p^p \\ &\geq \lambda \cdot \alpha(\tilde{\Delta}^{(j)}), \text{ Lemma 4.6,} \\ &\geq \lambda \cdot \Omega_p \left(\kappa^{-\frac{1}{\min\{p,2\}-1}} \right) \alpha(\Delta^*) \\ &= \Omega_p \left(\kappa^{-\frac{1}{\min\{p,2\}-1}} \right) \alpha(\Delta^*). \end{aligned}$$

From Lemma 4.6 we know,

$$\alpha(\lambda\tilde{\Delta}^{(i)}) \geq \|\mathbf{x}\|_p^p - \|\mathbf{x} - \lambda\tilde{\Delta}^{(i)}\|_p^p.$$

We thus have $\alpha(\lambda\tilde{\Delta}^{(i)}) \geq \Omega_p \left(\kappa^{-\frac{1}{\min\{p,2\}-1}} \right) \text{OPT}$, implying $\lambda \cdot \tilde{\Delta}^{(i)}$ is also a $\Omega_p \left(\kappa^{-\frac{1}{\min\{p,2\}-1}} \right)$ approximate solution as required. \square

It remains to solve problems of the form (2) up to a κ -approximation. Recall that these problems look like,

$$\begin{aligned} \min_{\Delta} \quad &\gamma_p(\mathbf{t}, \Delta) \\ &\hat{\mathbf{A}}\Delta = \mathbf{d}, \end{aligned}$$

and satisfy $\text{OPT} \leq 1$, and $m^{-1/p} \leq \mathbf{t}_j \leq 1, \forall j$.

5.2 Oracle

Our approach follows the format of the approximate max-flow algorithm by Christiano et al. [7]. Specifically, we use a variant of multiplicative weights update to converge to a solution with small $\gamma_p(\mathbf{t}, \Delta)$. The multiplicative weights update scheme repeatedly updates a set of weights \mathbf{w} using partial, local solutions computed based on these weights. The Christiano et al. algorithm can be viewed as picking these weights from the gradients of the soft-max function on flows. We will adapt this routine by showing that \mathbf{w} 's chosen from the gradient of $\gamma_p(\mathbf{t}, \Delta)$ also suffices for approximately minimizing the problem stated in 2.

The subroutine that this algorithm passes the \mathbf{w} onto is commonly referred to as an oracle. An oracle needs to compute a solution with both small dot-product against \mathbf{w} , and small *width*, which is defined as the maximum value of an entry. In such an oracle, the dot product condition is the hard constraint, in that the final approximation factor of the solution is directly related to the value of these dot products. The width, on the other hand, only affects the overall iteration count/running time, and can even be manipulated/improved algorithmically. Therefore we first need to define and show a good upper bound on the objective of the optimization problem solved within the oracle.

Formally, our oracle subroutine Algorithm 3 takes as input some affine constraints and vector of weights \mathbf{w} . It first computes a vector of non-negative weights \mathbf{r} , and then returns a minimizer

to the following optimization problem

$$\begin{aligned} \Delta &= \arg \min_{\Delta \in \mathbb{R}^m} \sum_e \mathbf{r}_e \Delta_e^2 \\ \text{s.t.} \quad & \widehat{\mathbf{A}} \Delta = \mathbf{d}. \end{aligned} \tag{4}$$

Appendix C contains an algorithm that solves such problems efficiently.

Algorithm 3 Oracle

- 1: **procedure** ORACLE($\widehat{\mathbf{A}}, \mathbf{d}, \mathbf{w}, \mathbf{t}$)
- 2: $\mathbf{r}_e \leftarrow \left(m^{1/p} \mathbf{t}_e\right)^{p-2} + \mathbf{w}_e^{p-2}$
- 3: Compute the Δ using resistances \mathbf{r}_e that satisfies solves the following optimization problem

$$\begin{aligned} \Delta &= \arg \min_{\Delta'} \sum_e \mathbf{r}_e \Delta_e'^2 \\ \text{s.t.} \quad & \widehat{\mathbf{A}} \Delta' = \mathbf{d} \end{aligned}$$

- 4: **return** Δ
-

Let us now look at some properties of the solution returned by the oracle. Note that the objective of our problem (2) is at most 1. This implies that we have Δ such that

- $\sum_e (\Delta_e^*)^2 \mathbf{t}_e^{p-2} \leq 1$,
- $\sum_e |\Delta_e^*|^p \leq 1$, or $\|\Delta^*\|_p \leq 1$.

We next look at some relations on the weights and resistances. The following lemma is a simple application of Hölder's inequality. Its proof is given in Appendix A.3.

Lemma 5.6. *Let $p \geq 2$. For any set of weights \mathbf{w} on the edges, $\sum_e \mathbf{w}_e^{p-2} (\Delta_e^*)^2 \leq \|\mathbf{w}\|_p^{p-2}$.*

Lemma 5.7. *Let $p \geq 2$. For any \mathbf{w} , let Δ be the electrical flow computed with respect to resistances*

$$\mathbf{r}_e \stackrel{\text{def}}{=} \left(m^{1/p} \mathbf{t}_e\right)^{p-2} + \mathbf{w}_e^{p-2},$$

and demand vector \mathbf{d} .

Then the following hold,

1. $\sum_e \Delta_e^2 \leq \sum_e \mathbf{r}_e \Delta_e^2 \leq m^{\frac{p-2}{p}} + \|\mathbf{w}\|_p^{p-2}$,
2. $\sum_e |\Delta_e| \left| \gamma'(m^{1/p} \mathbf{t}_e, \mathbf{w}_e) \right| \leq p \left(\sum_e \gamma_p(m^{1/p} \mathbf{t}_e, \mathbf{w}_e) \right)^{\frac{p-1}{p}} + pm^{\frac{p-2}{2p}} \left(\sum_e \gamma_p(m^{1/p} \mathbf{t}_e, \mathbf{w}_e) \right)^{\frac{1}{2}}$.

Proof. Since Δ is the electrical flow,

$$\sum_e \mathbf{r}_e \Delta_e^2 \leq \sum_e \mathbf{r}_e (\Delta_e^*)^2.$$

We have,

$$\begin{aligned} \sum_e \mathbf{r}_e \Delta_e^2 &\leq \sum_e \mathbf{r}_e (\Delta_e^*)^2 = \sum_e (m^{1/p} \mathbf{t}_e)^{p-2} (\Delta_e^*)^2 + \sum_e \mathbf{w}_e^{p-2} (\Delta_e^*)^2 \\ &\leq m^{\frac{p-2}{p}} + \|\mathbf{w}\|_p^{p-2}, \text{ follows from lemma 5.6 and the fact } \sum_e (\Delta_e^*)^2 \mathbf{t}_e^{p-2} \leq 1. \end{aligned}$$

Finally, using $\mathbf{r}_e \geq (m^{1/p} \mathbf{t}_e)^{p-2} \geq 1$, we have $\sum_e \Delta_e^2 \leq \sum_e \mathbf{r}_e \Delta_e^2$, completing part 1.

Now we know that,

$$|\gamma'(m^{1/p} \mathbf{t}_e, \mathbf{w}_e)| = \begin{cases} p(m^{1/p} \mathbf{t}_e)^{p-2} \mathbf{w}_e & \text{if } |\mathbf{w}_e| \leq m^{1/p} \mathbf{t}_e, \\ p|\mathbf{w}_e|^{p-2} \mathbf{w}_e & \text{otherwise.} \end{cases}$$

Using Cauchy Schwarz's inequality,

$$\begin{aligned} \left(\sum_e |\Delta_e| |\gamma'(m^{1/p} \mathbf{t}_e, \mathbf{w}_e)| \right)^2 &= \left(\sum_e p |\Delta_e| |\max(m^{1/p} \mathbf{t}_e, \mathbf{w}_e)|^{p-2} \mathbf{w}_e^2 \right)^2 \\ &\leq p^2 \left(\sum_e \max(m^{1/p} \mathbf{t}_e, \mathbf{w}_e)^{p-2} \mathbf{w}_e^2 \right) \left(\sum_e \max(m^{1/p} \mathbf{t}_e, \mathbf{w}_e)^{p-2} \Delta_e^2 \right) \\ &\leq p^2 \gamma_p(m^{1/p} \mathbf{t}, \mathbf{w}) \sum_e \max(m^{1/p} \mathbf{t}_e, |\mathbf{w}_e|)^{p-2} \Delta_e^2 \end{aligned}$$

Combining the two cases we have,

$$\begin{aligned} \sum_e |\Delta_e| |\gamma'(m^{1/p} \mathbf{t}_e, \mathbf{w}_e)| &\leq p \sqrt{\gamma_p(m^{1/p} \mathbf{t}, \mathbf{w}) \sum_e \max(m^{1/p} \mathbf{t}_e, |\mathbf{w}_e|)^{p-2} \Delta_e^2} \\ &\leq p \sqrt{\gamma_p(m^{1/p} \mathbf{t}, \mathbf{w}) \sum_e \mathbf{r}_e \Delta_e^2} \\ &\leq p \sqrt{\gamma_p(m^{1/p} \mathbf{t}, \mathbf{w})} \sqrt{m^{\frac{p-2}{p}} + \|\mathbf{w}\|_p^{p-2}} \\ &\leq pm^{\frac{p-2}{2p}} \sqrt{\gamma_p(m^{1/p} \mathbf{t}, \mathbf{w})} + p \sqrt{\gamma_p(m^{1/p} \mathbf{t}, \mathbf{w})} \|\mathbf{w}\|_p^{\frac{p-2}{2}} \\ &= pm^{\frac{p-2}{2p}} \gamma_p(m^{1/p} \mathbf{t}, \mathbf{w})^{\frac{1}{2}} + p \gamma_p(m^{1/p} \mathbf{t}, \mathbf{w})^{\frac{p-1}{p}}, \end{aligned}$$

where the last line uses $\|x\|_p^p \leq \gamma_p(m^{1/p} \mathbf{t}, \mathbf{w})$ for any \mathbf{t} . □

5.3 The Algorithm

Next, we integrate this oracle into the overall algorithm that repeatedly adjusts the weights. As with the use of electrical flow oracles for approximate max-flow [7], the convergence of such a scheme depends on the maximum values in the Δ returned by the oracle. However, because the overall objective is now a p -norm, the exact term of importance is actually the p -norm of Δ . Up to this discrepancy, we follow the algorithmic template from [7] by making an update when $\|\Delta\|_p^p$ is small, and make progress via another potential function otherwise.

In the cases where we do not take the step due to entries with large values, we show significant increases in an additional potential function, namely the objective of the quadratic minimization problem inside the oracle (Algorithm 3). However, the less graduate update schemes related to p -norms makes it no longer sufficient to update only the weight corresponding to the entry with maximum value. Furthermore, there may be entries with large values, whose corresponding resistances are too large for us to afford increasing. We address this by a scheme where we update an entry only if its value is larger than some threshold ρ , and that its resistance is at most another threshold β . Specifically, we show that for an appropriate choice of β and ρ , such updates both do not change the primary potential function (related to $\gamma_p(t, \mathbf{x})$) by too much (in Lemma 5.10), and increases the secondary potential function (the objective of the quadratic minimization problem) significantly whenever $\|\Delta\|_p^p$ is large (in Lemma 5.13). Pseudocode of this scheme is in Algorithm 4.

Algorithm 4 Algorithm for the Scaled down Problem

```

1: procedure GAMMA-SOLVER( $\mathbf{A}'$ ,  $\mathbf{c}$ ,  $\mathbf{t}$ )
2:    $\mathbf{w}_e^{(0,0)} \leftarrow 0$ 
3:    $\mathbf{x} \leftarrow 0$ 
4:    $\rho \leftarrow \tilde{\Theta}_p \left( m^{\frac{(p^2-4p+2)}{p(3p-2)}} \right)$  ▷ width parameter
5:    $\beta \leftarrow \tilde{\Theta}_p \left( m^{\frac{p-2}{3p-2}} \right)$  ▷ resistance threshold
6:    $\alpha \leftarrow \tilde{\Theta}_p \left( m^{-\frac{p^2-5p+2}{p(3p-2)}} \right)$  ▷ step size
7:    $\tau \leftarrow \tilde{\Theta}_p \left( m^{\frac{(p-1)(p-2)}{(3p-2)}} \right)$  ▷  $\ell_p$  energy threshold
8:    $T \leftarrow \alpha^{-1} m^{1/p} = \tilde{\Theta}_p \left( m^{\frac{p-2}{3p-2}} \right)$ 
9:    $i \leftarrow 0, k \leftarrow 0$ 
10:  while  $i < T$  do
11:     $\Delta = \text{ORACLE}(\mathbf{A}', \mathbf{c}, \mathbf{w}^{(i,k)}, \mathbf{t})$ 
12:    if  $\|\Delta\|_p^p \leq \tau$  then ▷ flow step
13:       $\mathbf{w}^{(i+1,k)} \leftarrow \mathbf{w}^{(i,k)} + \alpha |\Delta|$ 
14:       $\mathbf{x} \leftarrow \mathbf{x} + \alpha \Delta$ 
15:       $i \leftarrow i + 1$ 
16:    else ▷ width reduction step
17:      For all edges  $e$  with  $|\Delta_e| \geq \rho$  and  $\mathbf{r}_e \leq \beta$ 
18:         $\mathbf{w}_e^{(i,k+1)} \leftarrow 4^{\frac{1}{p-2}} \max(m^{1/p} \mathbf{t}_e, \mathbf{w}_e^{(i,k)})$ 
19:       $k \leftarrow k + 1$ 
20:  return  $m^{-\frac{1}{p}} \mathbf{x}$ 

```

Theorem 5.8. Let $p \geq 2$. Given a matrix $\hat{\mathbf{A}}$ and vectors \mathbf{x} and \mathbf{t} such that $\forall e, m^{-1/p} \leq \mathbf{t}_e \leq 1$, Algorithm 4 uses $O_p \left(m^{\frac{p-2}{(3p-2)}} \left(\log \left(\frac{m \|\mathbf{d}\|_2^2}{\|\hat{\mathbf{A}}\|_2^2} \right) \right)^{\frac{p}{3p-2}} \right)$ calls to the oracle and returns a vector \mathbf{x} such

that $\widehat{\mathbf{A}}\mathbf{x} = \mathbf{d}$, and $\gamma_p(\mathbf{t}, \mathbf{x}) = O_p(1)$.

Analysis of Potentials.

We define the following potential function for the analysis of our algorithm.

Definition 5.9. Let Φ be the potential function defined as

$$\Phi(\mathbf{w}^{(i)}) \stackrel{\text{def}}{=} \gamma_p(m^{1/p}\mathbf{t}, \mathbf{w}^{(i)}).$$

Initially, since we start with $\mathbf{w}^{(0)} = 0$, we have $\Phi(\mathbf{w}^{(0)}) = 0$. Observe that in the algorithm, we update the potentials in both the flow step and the width reduction step whereas we update the solution only in the flow step. It is easy to see that we always have $\mathbf{w}^{(i,k)} \geq \left\lfloor \mathbf{x}^{(i,k)} \right\rfloor$.

We next bound the potential. In addition, we track the energy of the electrical flow in the network with resistances \mathbf{r} . Let $\Psi(\mathbf{r})$ denote the minimum of routing \mathbf{d} with resistances \mathbf{r} :

$$\Psi(\mathbf{r}) \stackrel{\text{def}}{=} \min_{\Delta: A'\Delta = \mathbf{d}} \sum_e r_e \Delta_e^2. \quad (5)$$

Note that this energy is equal to the energy calculated using the Δ obtained in the solution of (4).

Notation. We overload notation for $\Psi(i, k)$ to denote $\Psi(\mathbf{r}^{(i,k)})$.

Our proof of Theorem 5.8 will be based two main parts:

1. Provided the total number of width reduction steps, K , is not too big, then $\Phi(T, K)$ is small. This in turn upper bounds cost of the approximate solution $m^{-1/p}\mathbf{x}$.
2. Showing that K cannot be too big, because each width reduction step cause large growth in $\Psi(\cdot)$, while we can bound the total growth in $\Psi(\cdot)$ by relating it to $\Phi(\cdot)$.

We start by observing that when we when increase the weight \mathbf{w}_e of an edge during a width reduction step, this has the effect of at least doubling the resistance \mathbf{r}_e . Recall,

$$\mathbf{r}_e^{(i,k)} \stackrel{\text{def}}{=} \left(m^{1/p}\mathbf{t}_e\right)^{p-2} + \left(\mathbf{w}_e^{(i,k)}\right)^{p-2}.$$

Now,

$$\frac{\mathbf{r}_e^{(i,k+1)}}{\mathbf{r}_e^{(i,k)}} = \frac{\left(m^{1/p}\mathbf{t}_e\right)^{p-2} + \left(\mathbf{w}_e^{(i,k+1)}\right)^{p-2}}{\left(m^{1/p}\mathbf{t}_e\right)^{p-2} + \left(\mathbf{w}_e^{(i,k)}\right)^{p-2}} = \frac{\left(m^{1/p}\mathbf{t}_e\right)^{p-2} + 4 \max\left\{m^{1/p}\mathbf{t}_e, \mathbf{w}_e^{(i,k)}\right\}^{p-2}}{\left(m^{1/p}\mathbf{t}_e\right)^{p-2} + \left(\mathbf{w}_e^{(i,k)}\right)^{p-2}} \geq 2. \quad (6)$$

Meanwhile, the resistance does not grow by a factor larger than 4:

$$\frac{\mathbf{r}_e^{(i,k+1)}}{\mathbf{r}_e^{(i,k)}} = \frac{\left(m^{1/p}\mathbf{t}_e\right)^{p-2} + 4 \max\left\{m^{1/p}\mathbf{t}_e, \mathbf{w}_e^{(i,k)}\right\}^{p-2}}{\left(m^{1/p}\mathbf{t}_e\right)^{p-2} + \left(\mathbf{w}_e^{(i,k)}\right)^{p-2}} \leq 4. \quad (7)$$

We next show through the following lemma that the Φ potential does not increase too rapidly. The proof is through induction and can be found in Appendix B .

Lemma 5.10. *After i flow steps, and k width-reduction steps, provided*

1. $\alpha^p \tau \leq \alpha m^{\frac{p-1}{p}}$, (controls Φ growth in flow-steps)
2. $k \leq \rho^2 m^{2/p} \beta^{-\frac{2}{p-2}}$, (acceptable number of width-reduction steps)

the potential Φ is bounded as follows:

$$\Phi(i, k) \leq \left(p^2 2^p \alpha i + m^{1/p} \right)^p \exp \left(O_p(1) \frac{k}{\rho^2 m^{2/p} \beta^{-\frac{2}{p-2}}} \right).$$

We next wish to prove that in each width-reduction step, the electrical energy $\Psi(\cdot)$ goes up significantly. For this, we will use the following Lemma which is proven in Appendix D. It generalizes Lemma 2.6 of [7] to arbitrary weighted ℓ_2 regression problems, and directly measures the change in terms of the electrical energy of the entries modified.

Lemma 5.11. *Assuming the program (4) is feasible, let Δ be an be a solution to the optimization problem (4) with weights \mathbf{r} . Suppose we increase the resistance on each entry to get \mathbf{r}' . Then,*

$$\Psi(\mathbf{r}') \geq \exp \left(\frac{\sum_e \min \left\{ 1, \frac{r'_e - r_e}{r_e} \right\} r_e \Delta_e^2}{2\Psi(\mathbf{r})} \right) \Psi(\mathbf{r}).$$

This statement also implies the form of the lemma that concerns increasing the resistances on a set of entries uniformly [7, Lemma 2.6].

The next lemma gives a lower bound on the energy in iteration 0, i.e., when we start, and an upper bound on the energy at each step.

Lemma 5.12. *Initially, we have,*

$$\Psi(\mathbf{r}^{(0,0)}) \geq \frac{\|\mathbf{d}\|_2^2}{\|\mathbf{A}\|_2^2},$$

where $\|\mathbf{A}\|$ is the operator norm, or maximum singular value of \mathbf{A} . Let us call this ratio L . Moreover, at any step i , we have,

$$\Psi(\mathbf{r}^{(i,k)}) \leq m^{\frac{p-2}{p}} + \Phi(i, k)^{\frac{p-2}{p}}.$$

Proof. For the lower bound in the initial state, recall that we scale the problem such that $\text{OPT} = 1$, and $\mathbf{t}_e \geq m^{1/p}$. Initially we have, $\mathbf{r}_e^{(0,0)} = (m^{1/p} \mathbf{t}_e)^{p-2} \geq 1$. This means for any solution Δ , we have

$$\sum_e \mathbf{r}_e^{(0,0)} \Delta_e^2 \geq \|\Delta\|_2^2.$$

On the other hand, because

$$\|\mathbf{A}\Delta\|_2 \leq \|\mathbf{A}\|_2 \|\Delta\|_2,$$

we get

$$\|\Delta\|_2 \geq \frac{\|\mathbf{d}\|_2}{\|\mathbf{A}\|_2},$$

upon which squaring gives the lower bound on $\Psi(\mathbf{r}^{(0,0)})$.

For the upper bound, Lemma 5.7 implies that

$$\Psi(\mathbf{r}^{(i,k)}) \leq m^{\frac{p-2}{p}} + \|\mathbf{w}\|_p^{p-2} \leq m^{\frac{p-2}{p}} + \Phi(i,k)^{\frac{p-2}{p}}.$$

□

The next Lemma says that the three assumptions (stated in the statement of the Lemma) can be used to ensure that the potential $\Psi(\cdot)$ grows quickly with each width reduction step, and that flow steps do not cause the potential to shrink.

Lemma 5.13. *Suppose at step (i, k) we have $\|\Delta\|_p > \tau$ so that we perform a width reduction step (line 18). If*

1. $\Phi(i, k) \leq O_p(1)m$,
2. $\tau^{2/p} \geq 2\Omega_p(1)m^{\frac{p-2}{\beta}}$, and
3. $\frac{\tau}{10} \geq \rho^{p-2}m^{\frac{p-2}{p}}$.

Then

$$\Psi(i, k+1) \geq \Psi(i, k) \left(1 + \Omega_p(1) \frac{\tau^{2/p}}{m^{\frac{p-2}{p}}} \right).$$

Furthermore, if at (i, k) we have $\|\Delta\|_p \leq \tau$ so that we perform a flow step, then

$$\Psi(i+1, k) \geq \Psi(i, k).$$

Proof. It will be helpful for our analysis to split the index set into three disjoint parts:

- $S = \{e : |\Delta_e| \leq \rho\}$
- $H = \{e : |\Delta_e| > \rho \text{ and } \mathbf{r}_e \leq \beta\}$
- $B = \{e : |\Delta_e| > \rho \text{ and } \mathbf{r}_e > \beta\}$.

Firstly, we note

$$\sum_{e \in S} |\Delta_e|^p \leq \rho^{p-2} \sum_{e \in S} |\Delta_e|^2 \leq \rho^{p-2} \sum_{e \in S} \mathbf{r}_e |\Delta_e|^p \leq \rho^{p-2} O_p(1) m^{(p-2)/p}.$$

hence, using Assumption 3

$$\sum_{e \in H \cup B} |\Delta_e|^p \geq \sum_e \Delta_e^p - \sum_{e \in S} |\Delta_e|^p \geq \tau - \rho^{p-2} m^{\frac{p-2}{p}} \geq 9\tau.$$

This means,

$$\sum_{e \in H \cup B} \Delta_e^2 \geq \left(\sum_{e \in H \cup B} \Delta_e^p \right)^{p/2} \geq \Omega_p(1) \tau^{2/p}.$$

Secondly we note that, using Assumption (1) and Lemma 5.7, we have

$$\sum_{e \in B} \Delta_e^2 \leq \beta^{-1} \sum_{e \in B} \mathbf{r}_e \Delta_e^2 \leq \beta^{-1} O_p(1) m^{\frac{p-2}{p}}.$$

So then, using Assumption 2,

$$\sum_{e \in H} \Delta_e^2 = \sum_{e \in H \cup B} \Delta_e^2 - \sum_{e \in B} \Delta_e^2 \geq \Omega_p(1) \tau^{2/p} - \beta^{-1} O_p(1) m^{\frac{p-2}{p}} \geq \Omega_p(1) \tau^{2/p}.$$

As $\mathbf{r}_e \geq 1$, this implies $\sum_{e \in H} \mathbf{r}_e \Delta_e^2 \geq \Omega_p(1) \tau^{2/p}$.

From Lemma 5.12 and Assumption 1 we have

$$\Psi(i, k) \leq O_p(1) m^{(p-2)/p}.$$

So then, combining our last two observations, and applying Lemma 5.11, we get

$$\Psi(i, k+1) \geq \Psi(i, k) \left(1 + \Omega_p(1) \frac{\tau^{2/p}}{m^{\frac{p-2}{p}}} \right).$$

Finally, for the ‘‘flow step’’ case, by applying Lemma 5.11 with H as the whole set of indices, $\delta = 1$ and $\gamma = 1$, we get that as the resistances only increase,

$$\Psi(i+1, k) \geq \Psi(i, k).$$

□

We are now ready to prove Theorem 5.8.

Proof of Theorem 5.8

Proof. We first observe that our parameter choices in the Algorithm 4 satisfy Assumption 1 of Lemma 5.10, namely, we can choose the parameters α and τ s.t.

- $\alpha \leftarrow \Theta_p \left(m^{-\frac{p^2-5p+2}{p(3p-2)}} \left(\log \left(\frac{m}{L} \right) \right)^{\frac{-p}{3p-2}} \right),$
- $\tau \leftarrow \Theta_p \left(m^{\frac{(p-1)(p-2)}{(3p-2)}} \left(\log \left(\frac{m}{L} \right) \right)^{\frac{p(p-1)}{3p-2}} \right),$

while ensuring $\alpha^p \tau \leq \alpha m^{\frac{p-2}{p}}$. This means by Lemma 5.10, that if the Algorithm completes after taking $T = \alpha^{-1} m^{1/p}$ flow steps and $K \leq \Omega_p(1) \rho^2 m^{2/p} \beta^{-\frac{2}{p-2}}$, when it returns, we have

$$\Phi(T, K) \leq m \left(p^2 2^p + 1 \right)^p e^1 \leq O_p(1) m,$$

This means that the algorithm returns $m^{-\frac{1}{p}} \mathbf{x}$ with

$$\gamma(\mathbf{t}, m^{-\frac{1}{p}} \mathbf{x}) = \frac{1}{m} \gamma(m^{\frac{1}{p}} \mathbf{t}, \mathbf{x}) \leq \frac{1}{m} \gamma(m^{\frac{1}{p}} \mathbf{t}, \mathbf{w}^{(T, K)}) = \frac{1}{m} \Phi(T, K) \leq O_p(1).$$

Note the only alternative is that the algorithm takes more than $\Omega_p(1)\rho^2 m^{2/p}\beta^{-\frac{2}{p-2}}$ width reduction steps (and possibly infinitely many such steps, hence never terminating).

We will now show this cannot happen, by deriving a contradiction from the assumption that the algorithm takes a width reduction step starting from step (i, k) where $i < T$ and $k = \rho^2 m^{2/p}\beta^{-\frac{2}{p-2}}$.

Since the conditions for Lemma 5.10 hold for all preceding steps, we must have $\Phi(i, k) \leq O_p(1)m$.

Additionally, we note that our parameter choice of $\beta = \Theta_p\left(m^{\frac{p-2}{3p-2}}\left(\log\left(\frac{m}{L}\right)\right)^{-\frac{2(p-1)}{3p-2}}\right)$ and $\rho = \Theta_p\left(m^{\frac{(p^2-4p+2)}{p(3p-2)}}\left(\log\left(\frac{m}{L}\right)\right)^{\frac{p(p-1)}{(p-2)(3p-2)}}\right)$ along with our choice of τ (see above), ensures that

$$\tau^{2/p} \geq 2\Omega_p(1)\frac{m^{\frac{p-2}{p}}}{\beta} \quad \text{and} \quad \frac{\tau}{10} \geq \rho^{p-2}m^{\frac{p-2}{p}}.$$

This means that at every step (j, l) preceding the current step, the conditions of Lemma 5.13 are satisfied, so we can prove by a simple induction that

$$\Psi(i, k) \geq \Psi(0, 0) \left(1 + \Omega_p(1)\frac{\tau^{2/p}}{m^{\frac{p-2}{p}}}\right)^k > \Psi(0, 0) \exp\left(\Omega_p(1)\frac{\tau^{2/p}}{m^{\frac{p-2}{p}}}k\right).$$

Since our parameter choices ensure $\Omega_p(1)\frac{\tau^{2/p}}{m^{\frac{p-2}{p}}}k > \Theta_p\left(\frac{m}{L}\right)$ this means

$$\Psi(i, k) > \Psi(0, 0) \cdot \Theta_p\left(\frac{m}{L}\right).$$

But this contradicts Lemma 5.12, since this Lemma, combined with $\Phi(i, k) \leq O_p(1)m$ gives

$$\frac{\Psi(i, k)}{\Psi(0, 0)} \leq O_p\left(m^{\frac{p-2}{p}}\right).$$

From this contradiction, we conclude that we never have more than $K = \Omega_p(1)\rho^2 m^{2/p}\beta^{-\frac{2}{p-2}}$ width reduction steps.

Now we observe that the total number of oracle calls in the algorithm is bounded by

$$T + K \leq \Theta_p\left(m^{\frac{p-2}{3p-2}}\left(\log\left(\frac{m}{L}\right)\right)^{\frac{p}{3p-2}}\right).$$

□

This concludes the analysis of our algorithm.

5.4 Proof of Theorem 5.1

Proof. Theorem 5.8 implies that we can solve Program (2) using Algorithm 4 to get an $O_p(1)$ -approximate solution in $\tilde{O}_p\left(m^{\frac{p-2}{3p-2}}\right)$ calls to the Oracle. Implementing the Oracle requires solving

a linear system, and hence can be implemented in $O(m+n)^\omega$ time where ω is the matrix multiplication constant (see the Appendix for a proof). Thus, we can find an $O_p(1)$ -approximate solution to (2) in total time

$$\tilde{O}_p \left((m+n)^{\omega + \frac{p-2}{3p-2}} \right).$$

Now, Theorem 5.2 implies that we can find an $O_p(1)$ -approximate solution to the residual problem (1) in total time,

$$\tilde{O}_p \left((m+n)^{\omega + \frac{p-2}{3p-2}} \log 1/\varepsilon \right).$$

Finally using Theorem 4.1 we can conclude that we have an ε -approximate solution to (*) in $\tilde{O}_p \left(\log \frac{1}{\varepsilon} \right)$ calls to a $O_p(1)$ -approximate solver to the residual problem (1). This gives us a total running time of,

$$\tilde{O}_p \left((m+n)^{\omega + \frac{p-2}{3p-2}} \log^2 1/\varepsilon \right).$$

□

We now have a complete algorithm for the p -norm regression problem that gives an ε -approximate solution.

6 Speedups for General Matrices via. Inverse Maintenance

If \mathbf{A} is an explicitly given, $m \times n$, matrix, we need to solve the quadratic minimization problem at each step. This can be solved via a linear systems solve in the matrix

$$\mathbf{A}^\top \mathbf{Diag}(\mathbf{r})^{-1} \mathbf{A}.$$

which takes $O((m+n)^\omega)$, where ω is the matrix multiplication constant. This directly gives a total running time cost of $\tilde{O}_p(m^{\frac{p-2}{3p-2}}(m+n)^\omega \log(1/\varepsilon))$, which for large values of p , along with the assumption of $\omega > 2.37$, exceeds 2.70.

This is more than the running time of about $O(mn^{1.5})$ of algorithms based on inverse maintenance [25, 26, 45]. In this section we show that the MWU routine from Section 5 can also benefit from fast inverse maintenance. Our main result is:

Theorem 6.1. *If \mathbf{A} is an explicitly given, m -by- n matrix with polynomially bounded condition numbers, and $p \geq 2$ Algorithm 4 as given in Section 5.3 can be implemented to run in total time*

$$\tilde{O}_p \left((m+n)^{\max\left\{\omega, 2 + \frac{p-(10-4\omega)}{3p-2}\right\}} \right).$$

A few remarks about this running time: the term that dominates depends on the comparison between $2/3$ and $10 - 4\omega$, or after manipulation, the comparison between ω and $7/3$:

1. For the current best value of $\omega > 7/3$, the second term is at most ω , so the total running time is about $(m+n)^\omega$.
2. If $\omega = 2$, then this running time is simply $(m+n)^{\frac{p-2}{3p-2}}$: same as resolving the linear system at each step.

3. If $\omega \leq 7/3$, then the overhead in the exponent on the second term is at most

$$\frac{p-2/3}{3p-2} = 1/3,$$

and this value approaches $\frac{p-2}{3p-2}$ as $\omega \rightarrow 2$.

Our algorithm is based on gradually updating the \mathbf{r} vector. First, note that $\mathbf{w}_e^{(i)}$'s, and thus $\mathbf{r}_e^{(i)}$'s are monotonically increasing. Secondly, for the $\mathbf{r}^{(i)}$ that do not double, we can replace with the original version while forming a factor 2 preconditioner. Thus, we only need to update the $\mathbf{r}^{(i)}$ entries that have significant increases. This update can be encapsulated by the following result on computing low rank perturbations to a matrix, which is a direct consequence of rectangular matrix multiplication and Woodbury matrix formula.

Lemma 6.2. *Given an m -by- n matrix \mathbf{A} , along with vectors $\hat{\mathbf{r}}$ and $\tilde{\mathbf{r}}$ that differ in k entries, as well as the matrix $\hat{\mathbf{Z}} = (\mathbf{A}^\top \text{Diag}(\hat{\mathbf{r}})^{-1} \mathbf{A})^{-1}$, we can construct $(\mathbf{A}^\top \text{Diag}(\tilde{\mathbf{r}})^{-1} \mathbf{A})^{-1}$ in $O(k^{\omega-2}(m+n)^2)$ time.*

Proof. Let S denote the entries that differ in $\hat{\mathbf{r}}$ and $\tilde{\mathbf{r}}$. Then we have

$$\mathbf{A}^\top \text{Diag}(\tilde{\mathbf{r}})^{-1} \mathbf{A} = \mathbf{A}^\top \text{Diag}(\hat{\mathbf{r}})^{-1} \mathbf{A} + \mathbf{A}_{:,S}^\top \left(\text{Diag}(\tilde{\mathbf{r}}_S)^{-1} - \text{Diag}(\hat{\mathbf{r}}_S)^{-1} \right) \mathbf{A}_{S,:}.$$

This is a low rank perturbation, so by Woodbury matrix identity we get:

$$\left(\mathbf{A}^\top \text{Diag}(\tilde{\mathbf{r}})^{-1} \mathbf{A} \right)^{-1} = \hat{\mathbf{Z}} - \hat{\mathbf{Z}} \mathbf{A}_{:,S}^\top \left(\left(\text{Diag}(\tilde{\mathbf{r}}_S)^{-1} - \text{Diag}(\hat{\mathbf{r}}_S)^{-1} \right)^{-1} + \mathbf{A}_{S,:} \hat{\mathbf{Z}} \mathbf{A}_{:,S}^\top \right)^{-1} \mathbf{A}_{S,:} \hat{\mathbf{Z}},$$

where we use $\hat{\mathbf{Z}}^\top = \hat{\mathbf{Z}}$ because $\mathbf{A}^\top \text{Diag}(\hat{\mathbf{r}})^{-1} \mathbf{A}$ is a symmetric matrix. To explicitly compute this matrix, we need to:

1. compute the matrix $\mathbf{A}_{S,:} \hat{\mathbf{Z}}$,
2. compute $\mathbf{A}_{:,S} \hat{\mathbf{Z}} \mathbf{A}_{:,S}^\top$
3. invert the middle term.

This cost is dominated by the first term, which can be viewed as multiplying $\lceil n/k \rceil$ pairs of $k \times n$ and $n \times k$ matrices. Each such multiplication takes time $k^{\omega-1}n$, for a total cost of $O(k^{\omega-2}n^2)$. The other terms all involve matrices with dimension at most $k \times n$, and are thus lower order terms. \square

Note that the running time of Lemma 6.2 favours ‘batching’ a large number of modified edges to insert. To this end, we show that it suffices to have an inverse that only approximates some entries of $\mathbf{r}^{(i)}$. To do so, we first need to introduce our notions of approximations:

Definition 6.3. *We use $a \approx_c b$ for positive numbers a and b iff $c^{-1}a \leq b \leq c \cdot b$, and for vectors and for vectors \mathbf{a} and \mathbf{b} we use $\mathbf{a} \approx_c \mathbf{b}$ to denote $\mathbf{a}_i \approx_c \mathbf{b}_i$ entry-wise.*

Since we are only updating k resistances that have a constant factor increase and using a constant factor preconditioning for the others, we need the following result on preconditioned iterative methods for solving systems of linear equations.

Lemma 6.4. *If \mathbf{r} and $\hat{\mathbf{r}}$ are vectors such that $\mathbf{r} \approx_{\tilde{O}(1)} \hat{\mathbf{r}}$, and we're given the matrix $\hat{\mathbf{Z}} = (\mathbf{A}^\top \mathbf{Diag}(\hat{\mathbf{r}})^{-1} \mathbf{A})^{-1}$ explicitly, then we can solve a system of linear equations involving $\mathbf{A}^\top \mathbf{Diag}(\mathbf{r})^{-1} \mathbf{A}$ to $1/\text{poly}(n)$ accuracy in $\tilde{O}(n^2)$ time.*

As the resistances we provide to ORACLE are in the range $[1, O_p(m)]$, we get that each $\hat{\mathbf{r}}_e$ only needs to be updated $O(\log m)$ times, instead of after each iteration. However, it's insufficient to use this bound in the worst-case manner: if there are $m^{1/3}$ iterations, each of which doubles the resistances on $m^{2/3}$ edges, then the total cost as given by Lemma 6.2 becomes

$$(m+n)^{2+\frac{1}{3}+\frac{2}{3}\cdot(\omega-2)},$$

which is about $(m+n)^{2.58}$.

We get an even better bound by using our analysis from Section 5 to show that for iteration/edge combinations i and e , the (relative) update to $\mathbf{r}_e^{(i)}$ is small. Such small changes also imply that we can wait on such updates. For simplicity, suppose we only increment the resistances by factors of $\frac{1}{L}$, then it takes $\Theta(L)$ such increments until the edge's resistance has deviated by a constant factor. Furthermore, we can wait for another $\Theta(L)$ iterations before having to reflect this change in $\hat{\mathbf{r}}$: the total relative increases in these iterations is also at most $O(1)$. Formalizing this process leads to a lazy-update routine that tracks the increments of different sizes separately. Its Pseudocode is in Algorithms 5 and 6.

We will call the initialization routine INVERSEINIT at the first iteration, and subsequently call UPDATEINVERSE upon generating a new set of resistances in the call to Algorithm 3, ORACLE. This is in turn called from Line 11 of Algorithm 4. As a result, we will assume access to all variables of these routines. Furthermore, our routines keeps the following global variables:

1. $\hat{\mathbf{r}}$: resistances from the last time we updated each entry.
2. $counter(\eta)_e$: for each entry, track the number of times that it changed (relative to $\hat{\mathbf{r}}$) by a factor of about $2^{-\eta}$ since the previous update.
3. $\hat{\mathbf{Z}}$, an inverse of the matrix given by $\mathbf{A}^\top \mathbf{Diag}(\hat{\mathbf{r}})^{-1} \mathbf{A}$.

Algorithm 5 Inverse Maintenance Initialization

- 1: **procedure** INVERSEINIT
 - 2: Set $\hat{\mathbf{r}} \leftarrow \mathbf{r}^{(0)}$.
 - 3: Set $counter(\eta)_e \leftarrow 0$ for all $0 \leq \eta \leq \log(m)$ and e .
 - 4: Set $\mathbf{Z} \leftarrow (\mathbf{A}^\top \mathbf{Diag}(\hat{\mathbf{r}})^{-1} \mathbf{A})^{-1}$ by explicitly inverting the matrix.
-

Algorithm 6 Inverse Maintenance Procedure

```

1: procedure UPDATEINVERSE
2:   for all entries  $e$  do
3:     Find the least non-negative integer  $\eta$  such that
           
$$\frac{1}{2^\eta} \leq \frac{\mathbf{r}_e^{(i)} - \mathbf{r}_e^{(i-1)}}{\widehat{\mathbf{r}}_e}.$$

4:     Increment  $\text{counter}(\eta)_e$ .
5:      $E_{\text{changed}} \leftarrow \cup_{\eta: i \pmod{2^\eta} \equiv 0} \{e : \text{counter}(\eta)_e \geq 2^\eta\}$ 
6:      $\widetilde{\mathbf{r}} \leftarrow \widehat{\mathbf{r}}$ 
7:     for all  $e \in E_{\text{changed}}$  do
8:        $\widetilde{\mathbf{r}}_e \leftarrow \mathbf{r}_e^{(i)}$ .
9:       Set  $\text{counter}(\eta)_e \leftarrow 0$  for all  $\eta$ .
10:     $\widehat{\mathbf{Z}} \leftarrow \text{LOWRANKUPDATE}(\mathbf{A}, \mathbf{Z}, \widehat{\mathbf{r}}, \widetilde{\mathbf{r}})$ .
11:     $\widehat{\mathbf{r}} \leftarrow \widetilde{\mathbf{r}}$ .

```

We first verify that the maintained inverse is always a good preconditioner to the actual matrix, $\mathbf{A}^\top \mathbf{Diag}(\mathbf{r}^{(i)}) \mathbf{A}$.

Lemma 6.5. *After each call to UPDATEINVERSE, the vector $\widehat{\mathbf{r}}$ satisfies*

$$\widehat{\mathbf{r}} \approx_{\widetilde{O}(1)} \mathbf{r}^{(i)}.$$

Proof. First, observe that any change in resistance exceeding 1 is reflected immediately. Otherwise, every time we update $\text{counter}(j)_e$, \mathbf{r}_e can only increase additively by at most

$$2^{-j+1} \widehat{\mathbf{r}}_e.$$

Once $\text{counter}(j)_e$ exceeds 2^j , e will be added to E_{changed} after at most 2^j steps. So when we start from $\widehat{\mathbf{r}}_e$, e is added to E_{changed} after $\text{counter}(j)_e \leq 2^j + 2^j = 2^{j+1}$ iterations. The maximum possible increase in resistance due to the bucket j is,

$$2^{-j+1} \widehat{\mathbf{r}}_e \cdot 2^{j+1} = 4 \widehat{\mathbf{r}}_e.$$

Since there are only at most $m^{1/3}$ iterations, the contributions of buckets with $j > \log m$ are negligible. Now the change in resistance is influenced by all buckets j , each contributing at most $4 \widehat{\mathbf{r}}_e$ increase. The total change is at most $4 \widehat{\mathbf{r}}_e \log m$ since there are at most $\log m$ buckets. We therefore have

$$\widehat{\mathbf{r}}_e \leq \mathbf{r}_e^{(i)} \leq 5 \widehat{\mathbf{r}}_e \log m.$$

for every i . □

It remains to bound the number and sizes of calls made to Lemma 6.2. For this we define variables

$$k(\eta)^{(i)}$$

to denote the number of edges added to $E_{changed}$ at iteration i due to the value of $counter(\eta)_e$. Note that $k(\eta)^{(i)}$ is non-zero only if $i \equiv 0 \pmod{2^n}$, and

$$|E_{changed}^{(i)}| \leq \sum_{\eta} k(\eta)^{(i)}.$$

The following lemma gives a lower bound on the relative change of energy across one update of resistances.

Lemma 6.6. *Assuming the program (4) is feasible, let Δ be a solution to the optimization problem (4) with weights \mathbf{r} . Suppose we increase the resistance on each entry to \mathbf{r}' . Then,*

$$\frac{\Psi(\mathbf{r}') - \Psi(\mathbf{r})}{\Psi(\mathbf{r})} \geq \Omega_p \left(m^{-(p-2)/p} \sum_e (\Delta_e)^2 \min \left\{ 1, \frac{\mathbf{r}'_e - \mathbf{r}_e}{\mathbf{r}_e} \right\} \right).$$

Proof. Lemma 5.11 gives us that,

$$\begin{aligned} \frac{\Psi(\mathbf{r}')}{\Psi(\mathbf{r})} - 1 &\geq \exp \left(\frac{\sum_e \min \left\{ 1, \frac{\mathbf{r}'_e - \mathbf{r}_e}{\mathbf{r}_e} \right\} \mathbf{r}_e \Delta_e^2}{2\Psi(\mathbf{r})} \right) - 1 \\ &\geq \left(\frac{\sum_e \min \left\{ 1, \frac{\mathbf{r}'_e - \mathbf{r}_e}{\mathbf{r}_e} \right\} \mathbf{r}_e \Delta_e^2}{2\Psi(\mathbf{r})} \right) \end{aligned}$$

Since $\mathbf{r}_e \geq 1$ and $\Psi(\mathbf{r}) \leq O_p(m^{\frac{p-2}{p}})$,

$$\frac{\Psi(\mathbf{r}') - \Psi(\mathbf{r})}{\Psi(\mathbf{r})} \geq \Omega_p \left(m^{-(p-2)/p} \sum_e (\Delta_e)^2 \min \left\{ 1, \frac{\mathbf{r}'_e - \mathbf{r}_e}{\mathbf{r}_e} \right\} \right)$$

which gives our result. \square

We divide our analysis into 2 cases, when the relative change in resistance is at least 1 and when the relative change in resistance is at most 1. To begin with, let us first look at the following lemma that relates the change in weights to the relative change in resistance. The proof is in the Appendix.

Lemma 6.7. *Consider a flow step from Line 13 of Algorithm 4. We have*

$$\frac{\mathbf{r}_e^{(i+1)} - \mathbf{r}_e^{(i)}}{\mathbf{r}_e^{(i)}} \leq (1 + \alpha |\Delta_e|)^{p-2} - 1$$

where Δ is the ℓ_2 minimizer solution produced by the oracle.

Let us now see what happens when the relative change in resistance is at least 1.

Lemma 6.8. *Throughout the course of a run of Algorithm 4, the number of edges added to $E_{changed}$ due to relative resistance increase of at least 1,*

$$\sum_{1 \leq i \leq T} k(0)^{(i)} \leq \tilde{O}_P \left(m^{\frac{p+2}{3p-2}} \right).$$

Proof. From Lemma 6.6, we know that the relative change in energy over one iteration is at least,

$$\Omega_p \left(m^{-(p-2)/p} \sum_e (\Delta_e)^2 \min \left\{ 1, \frac{\mathbf{r}_e^{(i+1)} - \mathbf{r}_e^{(i)}}{\mathbf{r}_e^{(i)}} \right\} \right).$$

Over all iterations, the relative change in energy is at least,

$$\Omega_p \left(m^{-(p-2)/p} \sum_i \sum_e (\Delta_e)^2 \min \left\{ 1, \frac{\mathbf{r}_e^{(i+1)} - \mathbf{r}_e^{(i)}}{\mathbf{r}_e^{(i)}} \right\} \right)$$

which is upper bounded by $O(\log m)$. When iteration i is a width reduction step, the relative resistance change is always at least 1. In this case $|\Delta_e| \geq \rho$. When we have a flow step, Lemma 6.7 implies that when the relative change in resistance is at least 1 then,

$$\alpha |\Delta_e| \geq \Omega_p(1).$$

This gives, $|\Delta_e| \geq \Omega_p(\alpha^{-1})$. Using this bound on $|\Delta_e|$ is sufficient since $\rho > \Omega_p(\alpha^{-1})$ and both kinds of iterations are accounted for. The total relative change in energy can now be bounded.

$$\begin{aligned} & \Omega_p \left(m^{-(p-2)/p} \alpha^{-2} \sum_i \sum_e \mathbf{1}_{\left[\frac{\mathbf{r}_e^{(i+1)} - \mathbf{r}_e^{(i)}}{\mathbf{r}_e^{(i)}} \geq 1 \right]} \right) \leq \tilde{O}_p(1) \\ & \Leftrightarrow \Omega_p \left(m^{-(p-2)/p} \alpha^{-2} \sum_i k(0)^{(i)} \right) \leq \tilde{O}_p(1) \\ & \Leftrightarrow \sum_i k(0)^{(i)} \leq \tilde{O}_p(m^{(p-2)/p} \alpha^2). \end{aligned}$$

The Lemma follows by substituting $\alpha = \tilde{\Theta}_p \left(m^{-\frac{p^2-5p+2}{p(3p-2)}} \right)$ in the above equation. \square

Lemma 6.9. *Throughout the course of a run of Algorithm 4, the number of edges added to E_{changed} due to relative resistance increase between $2^{-\eta}$ and $2^{-\eta+1}$,*

$$\sum_{1 \leq i \leq T} k(\eta)^{(i)} \leq \begin{cases} 0 & \text{if } 2^\eta \geq T, \\ \tilde{O}_p \left(m^{\frac{p+2}{3p-2}} 2^{2\eta} \right) & \text{otherwise.} \end{cases}$$

Proof. From Lemma 6.6, the total relative change in energy is at least,

$$\Omega_p \left(m^{-(p-2)/p} \sum_i \sum_e (\Delta_e)^2 \left(\frac{\mathbf{r}_e^{(i+1)} - \mathbf{r}_e^{(i)}}{\mathbf{r}_e^{(i)}} \right) \right).$$

We know that $\frac{r_e^{(i+1)} - r_e^{(i)}}{r_e^{(i)}} \geq 2^{-\eta}$. Using Lemma 6.7, we have,

$$(1 + \alpha |\Delta_e|)^{p-2} - 1 \geq 2^{-\eta}.$$

We can bound $(1 + \alpha |\Delta_e|)^{p-2} - 1$ as,

$$(1 + \alpha |\Delta_e|)^{p-2} - 1 \leq \begin{cases} \alpha |\Delta_e| & \text{when } \alpha |\Delta_e| \leq 1 \text{ or } p-2 \leq 1 \\ O_p\left((\alpha |\Delta_e|)^{p-2}\right) & \text{otherwise.} \end{cases}$$

Now, in the second case, when $\alpha |\Delta_e| \geq 1$ and $p-2 > 1$,

$$(\alpha |\Delta_e|)^{p-2} \geq 2^{-\eta} \Rightarrow \alpha |\Delta_e| \geq \left(\frac{1}{2^\eta}\right)^{1/(p-2)} \geq 2^{-\eta}$$

For both cases we get,

$$\alpha |\Delta_e| \geq \Omega_p\left(2^{-\eta}\right).$$

Using the above bound and the fact that the total relative change in energy is at most $\tilde{O}_p(1)$, gives,

$$\begin{aligned} & \Omega_p \left(m^{-(p-2)/p} \sum_i \sum_e (\Delta_e)^2 \left(\frac{r_e^{(i+1)} - r_e^{(i)}}{r_e^{(i)}} \right) \right) \leq \tilde{O}_p(1) \\ \Rightarrow & \Omega_p \left(m^{-(p-2)/p} \sum_i \sum_e \left(\alpha^{-1} 2^{-\eta} \right)^2 \cdot \left(2^{-\eta} \mathbf{1}_{2^{-\eta+1} \geq \frac{r_e^{(i+1)} - r_e^{(i)}}{r_e^{(i)}} \geq 2^{-\eta}} \right) \right) \leq \tilde{O}_p(1) \\ \Rightarrow & \Omega_p \left(m^{-(p-2)/p} \alpha^{-2} 2^{-3\eta} \sum_i 2^\eta k(\eta)^{(i)} \right) \leq \tilde{O}_p(1) \\ \Rightarrow & \sum_i k(\eta)^{(i)} \leq \tilde{O}_p \left(m^{(p-2)/p} \alpha^2 2^{2\eta} \right) \end{aligned}$$

The Lemma follows substituting $\alpha = \tilde{\Theta}_p \left(m^{-\frac{p^2-5p+2}{p(3p-2)}} \right)$ in the above equation. \square

We can now use the concavity of $f(z) = z^{\omega-2}$ to upper bound the contribution of these terms.

Corollary 6.10. *Let $k(\eta)^{(i)}$ be as defined. Over all iterations we have,*

$$\sum_i \left(k(0)^{(i)} \right)^{\omega-2} \leq \tilde{O}_p \left(m^{\frac{p-(10-4\omega)}{3p-2}} \right)$$

and for every η ,

$$\sum_i^T \left(k(\eta)^{(i)} \right)^{\omega-2} \leq \begin{cases} 0 & \text{if } 2^\eta \geq T, \\ \tilde{O}_p \left(m^{\frac{p-2+4(\omega-2)}{3p-2}} \cdot 2^{\eta(3\omega-7)} \right) & \text{otherwise.} \end{cases}$$

Proof. Due to the concavity of the $\omega - 2 \approx 0.3727 < 1$ power, this total is maximized when it's equally distributed over all iterations. In the first sum, the number of terms is equal to the number of iterations, i.e., $\tilde{O}_p(m^{\frac{p-2}{3p-2}})$. In the second sum the number of terms is $\tilde{O}_p(m^{\frac{p-2}{3p-2}})2^{-\eta}$. Distributing the sum equally over the above numbers give,

$$\sum_i^T \left(k(0)^{(i)}\right)^{\omega-2} \leq \left(\tilde{O}_p\left(m^{\frac{p+2}{3p-2} - \frac{p-2}{3p-2}}\right)\right)^{\omega-2} \cdot \tilde{O}_p\left(m^{\frac{p-2}{3p-2}}\right) = \tilde{O}_p\left(m^{\frac{p-2+4(\omega-2)}{3p-2}}\right) \leq \tilde{O}_p\left(m^{\frac{p-(10-4\omega)}{3p-2}}\right)$$

and

$$\begin{aligned} \sum_i^T \left(k(\eta)^{(i)}\right)^{\omega-2} &\leq \tilde{O}\left(m^{\frac{p-2}{3p-2}}2^{-\eta}\right) \cdot \tilde{O}_p\left(\frac{m^{\frac{p+2}{3p-2}}2^{2\eta}}{m^{\frac{p-2}{3p-2}}2^{-\eta}}\right)^{\omega-2} \\ &= \tilde{O}_p\left(m^{\frac{p-2+4(\omega-2)}{3p-2}}2^{-\eta} \cdot 2^{3\eta(\omega-2)}\right) \\ &= \tilde{O}_p\left(m^{\frac{p-2+4(\omega-2)}{3p-2}}2^{\eta(3\omega-7)}\right). \end{aligned}$$

□

Proof. (of Theorem 6.1) By Lemma 6.5, the $\hat{\mathbf{r}}$ that the inverse being maintained corresponds to always satisfy $\hat{\mathbf{r}} \approx_{\tilde{O}(1)} \mathbf{r}^{(i)}$. So by the iterative linear systems solver method outlined in Lemma 6.4, we can implement each call to ORACLE (Section 5.2) in time $O((n+m)^2)$ in addition to the cost of performing inverse maintenance. This leads to a total cost of

$$\tilde{O}_p\left((n+m)^{2+\frac{p-2}{3p-2}}\right).$$

across the $T = \Theta_p(m^{\frac{p-2}{3p-2}})$ iterations.

The costs of inverse maintenance is dominated by the calls to the low-rank update procedure outlined in Lemma 6.2. Its total cost is bounded by

$$O\left(\sum_i \left|E_{\text{changed}}^{(i)}\right|^{\omega-2} (m+n)^2\right) = O\left((m+n)^2 \sum_i \left(\sum_{\eta} k(\eta)^{(i)}\right)^{\omega-2}\right).$$

Because there are only $O(\log m)$ values of η , and each $k(\eta)^{(i)}$ is non-negative, we can bound the total cost by:

$$\tilde{O}\left((m+n)^2 \sum_i \sum_{\eta} \left(k(\eta)^{(i)}\right)^{\omega-2}\right) \leq \tilde{O}_p\left((m+n)^2 \sum_{\eta: 2^\eta \leq T} m^{\frac{p-2+4(\omega-2)}{3p-2}} \cdot 2^{\eta(3\omega-7)}\right),$$

where the inequality follows from substituting in the result of Lemma 6.10. Depending on the sign of $3\omega - 7$, this sum is dominated either at $\eta = 0$ or $\eta = \log T$. Including both terms then gives

$$\tilde{O}_p\left((m+n)^{2+\frac{p-2+4(\omega-2)}{3p-2}} + (m+n)^{2+\frac{p-2+4(\omega-2)+(p-2)(3\omega-7)}{3p-2}}\right),$$

with the exponent on the trailing term simplifying to $\omega - 2$ to give,

$$\tilde{O}_p \left((m+n)^{2+\frac{p-(10-4\omega)}{3p-2}} + (m+n)^\omega \right).$$

□

7 Other Regression Formulations

In this section we discuss how various variants of ℓ_p -norm regression can be translated into our setting of

$$\begin{aligned} \min_{\mathbf{x}} \quad & \|\mathbf{x}\|_p \\ & \mathbf{A}\mathbf{x} = \mathbf{b}. \end{aligned} \tag{*}$$

As we will address a multitude of problems, we will make generic numerical assumptions to simplify the derivations.

1. $m \leq \text{poly}(n)$.
2. All entries in \mathbf{A} and \mathbf{b} are at most $\text{poly}(n)$. Note that this implies that the maximum singular value of \mathbf{A} , $\sigma_{\max}(\mathbf{A})$ and $\|\mathbf{b}\|_2$ are at most $\text{poly}(n)$.
3. $\|\mathbf{b}\|_2 \geq 1$.
4. The minimum non-zero singular value of \mathbf{A} , $\sigma_{\min}(\mathbf{A})$ is at least $1/\text{poly}(n)$.

Note that these conditions also imply bounds on the optimum value and the optimum solution \mathbf{x}^* : Specifically,

$$OPT \leq \|\mathbf{x}^*\|_2 \sqrt{m} \leq \|\mathbf{b}\|_2 \sigma_{\min}(\mathbf{A})^{-1} \leq \text{poly}(n),$$

and

$$\begin{aligned} OPT &\geq \|\mathbf{x}^*\|_2 \sqrt{m}^{-1} \geq \|\mathbf{b}\|_2 \sigma_{\max}(\mathbf{A})^{-1} \cdot m^{-1/2} \\ &\geq \frac{1}{\text{poly}(n)}. \end{aligned}$$

7.1 Affine transformations within the norm

Let \mathbf{C} be a matrix with the same assumptions as \mathbf{A} and \mathbf{d} have assumptions similar to \mathbf{b} . Suppose we are minimizing $\|\mathbf{C}\mathbf{x} - \mathbf{d}\|_p$ instead of $\|\mathbf{x}\|_p$, i.e.,

$$\begin{aligned} \min_{\mathbf{x}} \quad & \|\mathbf{C}\mathbf{x} - \mathbf{d}\|_p \\ & \mathbf{A}\mathbf{x} = \mathbf{b}. \end{aligned}$$

Note that this can be reduced to the following unconstrained problem,

$$\min_{\mathbf{x}} \|\mathbf{C}\mathbf{x} - \mathbf{d}\|_p.$$

To see this, first find the null space of \mathbf{A} , as well as a particular solution \mathbf{x}_0 that satisfies $\mathbf{A}\mathbf{x}_0 = \mathbf{b}$. Let the null space of \mathbf{A} be generated by the matrix \mathbf{V} . Then the space of solutions can be parameterized as

$$\mathbf{x} = \mathbf{x}_0 + \mathbf{V}\mathbf{y},$$

for some vector \mathbf{y} . Now our objective becomes,

$$\mathbf{C}\mathbf{V}\mathbf{y} + (\mathbf{C}\mathbf{x}_0 - \mathbf{d}),$$

which can be written as,

$$\min_{\mathbf{y}} \left\| \widehat{\mathbf{C}}\mathbf{y} - \widehat{\mathbf{d}} \right\|_p.$$

where $\widehat{\mathbf{C}} = \mathbf{C}\mathbf{V}$ and $\widehat{\mathbf{d}} = \mathbf{C}\mathbf{x}_0 - \mathbf{d}$. Observe that $\widehat{\mathbf{C}}\mathbf{y}$ spans the column space of $\widehat{\mathbf{C}}$. Decomposing $\widehat{\mathbf{d}}$ into a linear combination of an orthonormal basis we could combine the part which is in the span of $\widehat{\mathbf{C}}$ with $\widehat{\mathbf{C}}\mathbf{y}$. We can thus replace $\widehat{\mathbf{C}}$ in the objective with an orthonormal basis \mathbf{U} of its column space and replace $\widehat{\mathbf{d}}$ by \mathbf{g} , a vector orthogonal to all columns of \mathbf{U} . Then any vector

$$\mathbf{z} = \mathbf{U}\mathbf{y} - \mathbf{g},$$

can equivalently be described by the conditions

$$\begin{aligned} \mathbf{z}^\top \widehat{\mathbf{d}}_\perp &= \|\mathbf{g}\|_2^2, \\ \mathbf{z}^\top \mathbf{v} &= 0, \quad \forall \mathbf{v} \text{ s.t. } \mathbf{U}\mathbf{v} = 0, \mathbf{g}^\top \mathbf{v} = 0. \end{aligned}$$

For the last condition, it suffices to generate an orthonormal basis of the null space of \mathbf{U} . So the problem can be written as a linear constraint on \mathbf{z} instead.

7.2 $1 < p < 2$

In case $1 < p < 2$, we instead solve the dual problem:

$$\begin{aligned} \max_{\mathbf{y}} \quad & \mathbf{b}^\top \mathbf{y} \\ & \left\| \mathbf{A}^\top \mathbf{y} \right\|_q \leq 1, \end{aligned}$$

for $q = \frac{p}{p-1} > 2$. We can rescale the above problem to the equivalent q -norm ball-constrained projection problem,

$$\begin{aligned} \min_{\mathbf{y}} \quad & \left\| \mathbf{A}^\top \mathbf{y} \right\|_q \\ & \mathbf{b}^\top \mathbf{y} = 1, \end{aligned}$$

where the goal is to check whether the optimum is less than 1. This problem is covered by the problem introduced in Section 7.1 and can thus be solved to high accuracy in the desired time.

It remains is to transform a nearly-optimal solution \mathbf{y} of this q -norm ball-constrained projection problem to a nearly-optimal solution \mathbf{x} of the original subspace p -norm minimization problem. Since both of these problems' solutions are invariant under scalings to \mathbf{A} or \mathbf{b} , we may also assume that the optimum is at most 1.

Lemma 7.1. *If the optimum of*

$$\begin{aligned} \max \quad & \mathbf{b}^\top \mathbf{y} \\ & \left\| \mathbf{A}^\top \mathbf{y} \right\|_q \leq 1, \end{aligned}$$

is at most 1, and we have some \mathbf{y} such that

$$\begin{aligned} \mathbf{b}^\top \mathbf{y} &\geq 1 - \delta \\ \left\| \mathbf{A}^\top \mathbf{y} \right\|_q &= 1, \end{aligned}$$

then the gradient of $\left\| \mathbf{A}^\top \mathbf{y} \right\|_q$,

$$\nabla = \mathbf{A} \operatorname{sgn}(\mathbf{A}^\top \mathbf{y}) \left(\mathbf{A}^\top \mathbf{y} \right)^{q-1},$$

satisfies

$$\|\nabla - \mathbf{b}\|_2 \leq \delta \operatorname{poly}(n).$$

Proof. Let $\Delta = \nabla - \mathbf{b}$ and $p(n)$ be a polynomial such that $\gamma_q \left(\left| \mathbf{A}^\top \mathbf{y} \right|, \mathbf{A}^\top (\nabla - \mathbf{b}) \right) \leq p(n)$. By the assumption of \mathbf{A} and \mathbf{b} being $\operatorname{poly}(n)$ bounded, the above γ function is polynomially bounded. Let $q(n)$ be a polynomial in n such that, $q(n) \geq \sqrt{\frac{4p(n)}{\delta}}$. Suppose,

$$\|\nabla - \mathbf{b}\|_2^2 \geq \varepsilon > \delta q(n).$$

This gives us,

$$\Delta^\top \nabla \geq \Delta^\top \mathbf{b} + \varepsilon.$$

Now consider the solution

$$\hat{\mathbf{y}}(\theta) \leftarrow \mathbf{y} - \theta \Delta,$$

for step size $\theta = \varepsilon/2p(n)$. Lemma 4.5 and Lemma 3.3 gives

$$\begin{aligned} \left\| \mathbf{A}^\top \hat{\mathbf{y}} \right\|_q^q &\leq 1 - \theta \Delta^\top \nabla + \gamma_q \left(\left| \mathbf{A}^\top \mathbf{y} \right|, \theta \mathbf{A}^\top \Delta \right) \\ &\leq 1 - \theta \Delta^\top \nabla + \theta^2 \gamma_q \left(\left| \mathbf{A}^\top \mathbf{y} \right|, \mathbf{A}^\top \Delta \right) \\ &\leq 1 - \theta \Delta^\top \mathbf{b} - \theta \varepsilon + \theta^2 p(n) \\ &= 1 - \theta \Delta^\top \mathbf{b} - \frac{\theta \varepsilon}{2}. \end{aligned}$$

We can scale the solution $\hat{\mathbf{y}}$ up by a factor of $1/(1 - \theta \Delta^\top \mathbf{b} - \theta \varepsilon/2)$ to get a solution with objective value

$$\frac{1 - \delta - \theta \Delta^\top \mathbf{b}}{1 - \theta \Delta^\top \mathbf{b} - \theta \varepsilon/2}.$$

But by the assumption of 1 being the optimum, this cannot exceed 1, so we get

$$1 - \delta - \theta \Delta^\top \mathbf{b} \leq 1 - \theta \Delta^\top \mathbf{b} - \theta \varepsilon/2,$$

or

$$\theta\epsilon \leq 2\delta,$$

which combined with the choice of θ gives $\epsilon < \delta q(n)$ which is a contradiction. So we must have, $\|\nabla - \mathbf{b}\|_2 \leq \delta q(n) \leq \delta \text{poly}(n)$ \square

This means once $\delta \leq 1/\text{poly}(n)$, the solution created from the gradient

$$\hat{\mathbf{x}} \leftarrow \text{sgn}(\mathbf{A}^\top \mathbf{y}) \left(\mathbf{A}^\top \mathbf{y} \right)^{q-1},$$

satisfies

$$\|\mathbf{A}\hat{\mathbf{x}} - \mathbf{b}\|_2^2 \leq \delta \text{poly}(n).$$

Also, because $\sigma_{\min}(\mathbf{A}) \geq 1/\text{poly}(n)$, we can create a solution $\tilde{\mathbf{x}}$ from $\hat{\mathbf{x}}$ by doing a least squares projection on this difference. This gives:

$$\mathbf{A}\tilde{\mathbf{x}} = \mathbf{b},$$

and

$$\|\mathbf{x} - \hat{\mathbf{x}}\| \leq \delta \text{poly}(n) \sigma_{\min}(\mathbf{A}) \leq \delta \text{poly}(n).$$

Furthermore, note that because

$$\left(z^{q-1} \right)^p = z^{(q-1)(p-1)+q-1} = z^q,$$

we have

$$\|\nabla\|_p^p = \left\| \mathbf{A} \text{sgn}(\mathbf{A}^\top \mathbf{y}) \left(\mathbf{A}^\top \mathbf{y} \right)^{q-1} \right\|_p^p = \left\| \mathbf{A}^\top \mathbf{y} \right\|_q^q = 1,$$

so

$$\|\mathbf{x}\|_p \leq \|\nabla\|_p + \|\mathbf{x} - \nabla\|_p \leq 1 + \delta \text{poly}(n).$$

Thus, for sufficiently small δ , we can get high accuracy answer to the p -norm problem as well.

8 p -Norm Optimization on Graphs

In this section we discuss the performance of our algorithms on graphs. Here instead of invoking general linear algebraic routines, we instead invoke Laplacian solvers, which provide $1/\text{poly}(n)$ accuracy solutions in nearly-linear ($\tilde{O}(m)$) time [40].

Such matrices can be succinctly describe as

$$\mathbf{A}^\top \text{Diag}(\mathbf{r})^{-1} \mathbf{A}$$

where \mathbf{r} is the vector of resistances just as provided in the Oracle from Algorithm 3, but \mathbf{A} is the edge-vertex incidence matrix: with each row corresponding to an edge, each column corresponding to a vertex, and entries given by:

$$\mathbf{A}_{e,v} = \begin{cases} 1 & \text{if } v \text{ is the head of } e, \\ -1 & \text{if } v \text{ is the tail of } e, \\ 0 & \text{otherwise.} \end{cases}$$

Throughout this entire section, we will use \mathbf{A} to refer to the edge-vertex incidence matrix of a graph.

The main difficulty of reducing to Laplacian solvers is that we can no longer manipulate general matrices. Specifically, instead of directly working with the normal matrices as in Section 7.1, we need to implicitly track the subspaces, and optimize quadratics on them. As a result, we need to tailor such reductions towards the specific problems.

8.1 p -Norm Flows

This is closest to the general regression problem that we study:

$$\begin{aligned} \min_{\mathbf{x}} \quad & \|\mathbf{x}\|_p \\ & \mathbf{A}^\top \mathbf{x} = \mathbf{b} \end{aligned}$$

except with \mathbf{A} as an edge vertex incidence matrix.

When $p \geq 2$, the residual problem then has an extra condition of

$$\mathbf{g}^\top \mathbf{f} = \alpha,$$

which means we need to solve the problem of

$$\begin{aligned} \min_{\mathbf{x}} \quad & \sum_e r_e \mathbf{x}_e^2 \\ & \mathbf{A}^\top \mathbf{x} = \mathbf{b} \\ & \mathbf{g}^\top \mathbf{x} = \alpha \end{aligned}$$

which becomes a solve in the system of linear equations

$$[\mathbf{A}; \mathbf{x}]^\top \mathbf{Diag}(\mathbf{r})^{-1} [\mathbf{A}; \mathbf{x}] = \begin{bmatrix} \mathbf{A}^\top \mathbf{Diag}(\mathbf{r})^{-1} \mathbf{A} & \mathbf{A}^\top \mathbf{Diag}(\mathbf{r})^{-1} \mathbf{x} \\ \mathbf{x}^\top \mathbf{Diag}(\mathbf{r})^{-1} \mathbf{A} & \mathbf{x}^\top \mathbf{r}^{-1} \mathbf{x} \end{bmatrix}.$$

This matrix is a rank 3 perturbation to the graph Laplacian $\mathbf{A}^\top \mathbf{Diag}(\mathbf{r})^{-1} \mathbf{A}$, and can thus be solved in $\tilde{O}(m)$ time. A more detailed analysis of a generalization of this case can be found in Appendix B of [12].

When $1 < p < 2$, we invoke the dualization from Section 7.2 to obtain

$$\begin{aligned} \min_{\mathbf{y}} \quad & \|\mathbf{A}\mathbf{y}\|_q \\ & \mathbf{b}^\top \mathbf{y} = \alpha \end{aligned}$$

and if we retain the form of $\mathbf{A}^\top \mathbf{y}$, but transfer the gradient over to \mathbf{y} , the problem that we get is:

$$\begin{aligned} \min_{\mathbf{y}} \quad & \mathbf{y}^\top \mathbf{A}^\top \mathbf{Diag}(\mathbf{r})^{-1} \mathbf{A} \mathbf{y} \\ & \mathbf{b}^\top \mathbf{y} = \alpha \\ & \mathbf{g}^\top \mathbf{y} = \beta \end{aligned}$$

The two additional linear constraints can be removed by writing a variable of \mathbf{y} as a linear combination of the rest (as well as α). This then gives an unconstrained minimization problem on a subset of entries S ,

$$\min_{\mathbf{y}_S} \mathbf{y}_S^\top \mathbf{L}_{[S,S]} \mathbf{y}_S + \mathbf{c}^\top \mathbf{y}_S$$

where $\mathbf{L}_{[S,S]}$ is a minor of the Laplacian above, and this solution is obtained by solving for

$$\mathbf{y}_S \leftarrow \frac{1}{2} \mathbf{L}_{[S,S]}^\dagger \mathbf{c}.$$

8.2 Lipschitz Learning and Graph Labelling

This problem asks to label the vertices of a graph, with a set T fixed to the vector \mathbf{s} , while minimizing the p -norm difference between neighbours. It can be written as

$$\min_{\mathbf{x}: \mathbf{x}|_{T=\mathbf{s}}} \|\mathbf{A}\mathbf{x}\|_p^p$$

where \mathbf{A} is the edge-vertex incidence matrix.

In the case of $p \geq 2$, the residue problem becomes

$$\begin{aligned} \min_{\mathbf{x}: \mathbf{x}|_{T=\mathbf{s}}} \quad & \mathbf{x}^\top \mathbf{A}^\top \mathbf{D} \mathbf{A} \mathbf{x} \\ & \mathbf{g}^\top \mathbf{x} = \beta \end{aligned}$$

Here the gradient condition can be handled in the same way as with the voltage problem above: by fixing one additional entry of $V \setminus T$, and then solving an unconstrained quadratic minimization problem on the rest of the variables.

In the case of $1 < p < 2$, we first write down the problem as an unconstrained minimization problem on $V \setminus T$:

$$\min_{\mathbf{x}_{V \setminus T}} \left\| \mathbf{A}_{:, V \setminus T} \mathbf{x}_{V \setminus T} - \mathbf{A}_{:, T} \mathbf{s} \right\|_p.$$

Let $\mathbf{b} = \mathbf{A}_{:, T} \mathbf{s}$ and taking the dual gives:

$$\begin{aligned} \max \quad & \mathbf{b}^\top \mathbf{y} \\ & \|\mathbf{y}\|_q \leq 1 \\ & \left(\mathbf{A}^\top \right)_{V \setminus T, :} \mathbf{y} = 0 \end{aligned}$$

That is, solving for a small q -norm flow that maximizes the cost against \mathbf{b} , while also having 0 residues at the vertices not in T .

As $q > 2$, we can now invoke our main algorithm on \mathbf{y} . Upon binary search, and taking residual problems, we get ℓ_2 problems of the form

$$\begin{aligned} \min \quad & \sum_e r_e \mathbf{y}_e^2 \\ & \left(\mathbf{A}^\top \right)_{V \setminus T, :} \mathbf{y} = 0 \\ & \mathbf{b}^\top \mathbf{y} = \alpha \\ & \mathbf{g}^\top \mathbf{y} = \beta, \end{aligned}$$

which is solved by another low rank perturbation on a minor of the graph Laplacian.

References

- [1] Ittai Abraham, David Durfee, Ioannis Koutis, Sebastian Krinninger, and Richard Peng. “On Fully Dynamic Graph Sparsifiers”. In: *Symposium on Foundations of Computer Science (FOCS)*. Available at: <http://arxiv.org/abs/1604.02094>. 2016, pp. 335–344 (cit. on p. 5).
- [2] Zeyuan Allen-Zhu, Yuanzhi Li, Rafael Mendes de Oliveira, and Avi Wigderson. “Much Faster Algorithms for Matrix Scaling”. In: *58th IEEE Annual Symposium on Foundations of Computer Science, FOCS 2017, Berkeley, CA, USA, October 15-17, 2017*. Available at: <https://arxiv.org/abs/1704.02315>. 2017, pp. 890–901 (cit. on p. 4).
- [3] Sanjeev Arora, Elad Hazan, and Satyen Kale. “The Multiplicative Weights Update Method: a Meta-Algorithm and Applications.” In: *Theory of Computing* 8.1 (2012), pp. 121–164 (cit. on p. 5).
- [4] Owe Axelsson. *Iterative Solution Methods*. New York, NY: Cambridge University Press, 1994. ISBN: 0-521-44524-8 (cit. on p. 3).
- [5] Sébastien Bubeck, Michael B. Cohen, Yin Tat Lee, and Yuanzhi Li. “An Homotopy Method for Lp 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: ACM, 2018, pp. 1130–1137. ISBN: 978-1-4503-5559-9. DOI: 10.1145/3188745.3188776 (cit. on pp. 1, 2, 3, 4, 5, 46).
- [6] Hui Han Chin, Aleksander Madry, Gary L. Miller, and Richard Peng. “Runtime guarantees for regression problems”. In: *Proceedings of the 4th conference on Innovations in Theoretical Computer Science*. ITCS ’13. Available at <http://arxiv.org/abs/1110.1358>. Berkeley, California, USA: ACM, 2013, pp. 269–282. ISBN: 978-1-4503-1859-4 (cit. on pp. 1, 2, 4, 5).
- [7] Paul Christiano, Jonathan A. Kelner, Aleksander Madry, Daniel A. Spielman, and Shang-Hua Teng. “Electrical flows, laplacian systems, and faster approximation of maximum flow in undirected graphs”. In: *Proceedings of the 43rd annual ACM symposium on Theory of computing*. STOC ’11. Available at <http://arxiv.org/abs/1010.2921>. San Jose, California, USA: ACM, 2011, pp. 273–282. ISBN: 978-1-4503-0691-1. DOI: 10.1145/1993636.1993674 (cit. on pp. 2, 3, 4, 5, 12, 14, 17).
- [8] Michael B. Cohen, Yin Tat Lee, and Zhao Song. “Solving Linear Programs in the Current Matrix Multiplication Time”. In: *CoRR* abs/1810.07896 (2018). Available at: <http://arxiv.org/abs/1810.07896> (cit. on p. 3).
- [9] Michael B. Cohen, Aleksander Madry, Piotr Sankowski, and Adrian Vladu. “Negative-Weight Shortest Paths and Unit Capacity Minimum Cost Flow in $\tilde{O}(m^{10/7} \log W)$ Time (Extended Abstract)”. In: *Proceedings of the Twenty-Eighth Annual ACM-SIAM Symposium on Discrete Algorithms, SODA 2017, Barcelona, Spain, Hotel Porta Fira, January 16-19*. Available at: <https://arxiv.org/abs/1605.01717>. 2017, pp. 752–771 (cit. on p. 4).
- [10] Michael B. Cohen, Aleksander Madry, Dimitris Tsipras, and Adrian Vladu. “Matrix Scaling and Balancing via Box Constrained Newton’s Method and Interior Point Methods”. In: *58th IEEE Annual Symposium on Foundations of Computer Science, FOCS 2017, Berkeley, CA, USA, October 15-17, 2017*. Available at: <https://arxiv.org/abs/1704.02310>. 2017, pp. 902–913 (cit. on p. 4).

- [11] Michael B. Cohen and Richard Peng. “ ℓ_p Row Sampling by Lewis Weights”. In: *Proceedings of the Forty-Seventh Annual ACM on Symposium on Theory of Computing*. STOC ’15. Available at <http://arxiv.org/abs/1412.0588>. Portland, Oregon, USA: ACM, 2015, pp. 183–192. ISBN: 978-1-4503-3536-2. DOI: 10.1145/2746539.2746567 (cit. on pp. 3, 4).
- [12] Samuel I. Daitch and Daniel A. Spielman. “Faster approximate lossy generalized flow via interior point algorithms”. In: *Proceedings of the 40th annual ACM symposium on Theory of computing*. STOC ’08. Available at <http://arxiv.org/abs/0803.0988>. Victoria, British Columbia, Canada: ACM, 2008, pp. 451–460. ISBN: 978-1-60558-047-0. DOI: 10.1145/1374376.1374441 (cit. on p. 33).
- [13] Mohsen Ghaffari, Andreas Karrenbauer, Fabian Kuhn, Christoph Lenzen, and Boaz Patt-Shamir. “Near-Optimal Distributed Maximum Flow: Extended Abstract”. In: *Proceedings of the 2015 ACM Symposium on Principles of Distributed Computing, PODC 2015, Donostia-San Sebastián, Spain, July 21 - 23, 2015*. Available at: <https://arxiv.org/abs/1508.04747>. 2015, pp. 81–90 (cit. on p. 3).
- [14] Manoj Gupta and Richard Peng. “Fully Dynamic $1 + \epsilon$ -Approximate Matchings”. In: *54th Annual IEEE Symposium on Foundations of Computer Science, FOCS 2013, 26-29 October, 2013, Berkeley, CA, USA*. Available at <http://arxiv.org/abs/1304.0378>. 2013, pp. 548–557. DOI: 10.1109/FOCS.2013.65 (cit. on p. 5).
- [15] Van Enden Henson. “Multigrid methods nonlinear problems: an overview”. In: *Computational Imaging, Santa Clara, CA, USA, January 20, 2003*. 2003, pp. 36–48 (cit. on p. 3).
- [16] Carl T Kelley. *Iterative methods for optimization*. Vol. 18. Siam, 1999 (cit. on p. 3).
- [17] Jonathan A. Kelner, Yin Tat Lee, Lorenzo Orecchia, and Aaron Sidford. “An Almost-Linear-Time Algorithm for Approximate Max Flow in Undirected Graphs, and its Multicommodity Generalizations”. In: *Proceedings of the Twenty-Fifth Annual ACM-SIAM Symposium on Discrete Algorithms, SODA 2014, Portland, Oregon, USA, January 5-7, 2014*. Available at <http://arxiv.org/abs/1304.2338>. 2014, pp. 217–226 (cit. on p. 3).
- [18] Jonathan A. Kelner, Gary L. Miller, and Richard Peng. “Faster approximate multicommodity flow using quadratically coupled flows”. In: *Proceedings of the 44th symposium on Theory of Computing*. STOC ’12. Available at <http://arxiv.org/abs/1202.3367>. New York, New York, USA: ACM, 2012, pp. 1–18. ISBN: 978-1-4503-1245-5. DOI: 10.1145/2213977.2213979 (cit. on p. 4).
- [19] Dana A Knoll and David E Keyes. “Jacobian-free Newton–Krylov methods: a survey of approaches and applications”. In: *Journal of Computational Physics* 193.2 (2004), pp. 357–397 (cit. on p. 3).
- [20] R. Kyng, A. B. Rao, and S. Sachdeva. “Fast, Provable Algorithms for Isotonic Regression in all L_p -norms”. In: *NIPS*. 2015 (cit. on p. 4).
- [21] R. Kyng, A. B. Rao, S. Sachdeva, and D. A Spielman. “Algorithms for Lipschitz learning on graphs”. In: *COLT*. 2015 (cit. on pp. 3, 4).
- [22] François Le Gall. “Powers of tensors and fast matrix multiplication”. In: *Proceedings of the 39th international symposium on symbolic and algebraic computation*. Available at: <https://arxiv.org/abs/140> ACM. 2014, pp. 296–303 (cit. on p. 1).

- [23] Y. T. Lee, S. Rao, and N. Srivastava. “A new approach to computing maximum flows using electrical flows”. In: *STOC*. 2013 (cit. on p. 3).
- [24] Yin Tat Lee. “Uniform Sampling and Inverse Maintenance”. Talk at Michael Cohen Memorial Symposium, Available at: <https://simons.berkeley.edu/talks/welcome-and-birds-eye-view-michaels-work>. 2017 (cit. on p. 3).
- [25] Yin Tat Lee and Aaron Sidford. “Efficient Inverse Maintenance and Faster Algorithms for Linear Programming”. In: *IEEE 56th Annual Symposium on Foundations of Computer Science, FOCS 2015, Berkeley, CA, USA, 17-20 October, 2015*. Available at: <https://arxiv.org/abs/1503.01752>. 2015, pp. 230–249 (cit. on pp. 2, 4, 5, 21).
- [26] Yin Tat Lee and Aaron Sidford. “Path Finding Methods for Linear Programming: Solving Linear Programs in $\tilde{O}\sqrt{\text{rank}}$ Iterations and Faster Algorithms for Maximum Flow”. In: *Foundations of Computer Science (FOCS), 2014 IEEE 55th Annual Symposium on*. Available at <http://arxiv.org/abs/1312.6677> and <http://arxiv.org/abs/1312.6713>. IEEE. 2014, pp. 424–433 (cit. on pp. 1, 4, 5, 21).
- [27] Yin Tat Lee, Aaron Sidford, and Sam Chiu-wai Wong. “A Faster Cutting Plane Method and its Implications for Combinatorial and Convex Optimization”. In: *IEEE 56th Annual Symposium on Foundations of Computer Science, FOCS 2015, Berkeley, CA, USA, 17-20 October, 2015*. Available at: <https://arxiv.org/abs/1508.04874>. 2015, pp. 1049–1065 (cit. on p. 4).
- [28] A. Madry. “Navigating Central Path with Electrical Flows: From Flows to Matchings, and Back”. In: *FOCS*. 2013 (cit. on pp. 1, 3, 4).
- [29] Aleksander Madry. Private Communication. 2011 (cit. on p. 2).
- [30] Aleksander Madry. “Computing Maximum Flow with Augmenting Electrical Flows”. In: *IEEE 57th Annual Symposium on Foundations of Computer Science, FOCS 2016, 9-11 October 2016, Hyatt Regency, New Brunswick, New Jersey, USA*. Available at: <https://arxiv.org/abs/1608.06016>. 2016, pp. 593–602 (cit. on pp. 3, 4).
- [31] Y. Nesterov and A. Nemirovskii. *Interior-Point Polynomial Algorithms in Convex Programming*. Society for Industrial and Applied Mathematics, 1994. DOI: 10.1137/1.9781611970791. eprint: <https://epubs.siam.org/doi/pdf/10.1137/1.9781611970791> (cit. on p. 1).
- [32] Jorge Nocedal and Stephen J Wright. *Nonlinear Equations*. Springer, 2006 (cit. on p. 3).
- [33] Richard Peng. “Approximate undirected maximum flows in $O(m \text{ polylog}(n))$ time”. In: *Proceedings of the Twenty-Seventh Annual ACM-SIAM Symposium on Discrete Algorithms*. Available at <http://arxiv.org/abs/1411.7631>. SIAM. 2016, pp. 1862–1867 (cit. on pp. 1, 3).
- [34] Harald Racke, Chintan Shah, and Hanjo Taubig. “Computing Cut-Based Hierarchical Decompositions in Almost Linear Time”. In: *Proceedings of the 25th Annual ACM-SIAM Symposium on Discrete Algorithms*. SODA ‘14. 2014, pp. 227–238 (cit. on p. 3).
- [35] Y. Saad. *Iterative Methods for Sparse Linear Systems*. 2nd. Available at <http://www-users.cs.umn.edu/saad/> Philadelphia, PA, USA: Society for Industrial and Applied Mathematics, 2003. ISBN: 0898715342 (cit. on p. 3).
- [36] Jonah Sherman. “Area-convexity, l_∞ regularization, and undirected multicommodity flow”. In: *Proceedings of the 49th Annual ACM SIGACT Symposium on Theory of Computing, STOC 2017, Montreal, QC, Canada, June 19-23, 2017*. 2017, pp. 452–460 (cit. on pp. 1, 3).

- [37] Jonah Sherman. “Generalized Preconditioning and Undirected Minimum-cost Flow”. In: *Proceedings of the Twenty-Eighth Annual ACM-SIAM Symposium on Discrete Algorithms*. SODA ’17. Available at: <https://arxiv.org/abs/1606.07425>. Barcelona, Spain, 2017, pp. 772–780 (cit. on pp. 1, 3).
- [38] Jonah Sherman. “Nearly Maximum Flows in Nearly Linear Time”. In: *54th Annual IEEE Symposium on Foundations of Computer Science, FOCS 2013, 26-29 October, 2013, Berkeley, CA, USA*. Available at <http://arxiv.org/abs/1304.2077>. 2013, pp. 263–269 (cit. on p. 3).
- [39] Aaron Sidford. “Michael Cohen and Oblivious Routing”. Talk at Michael Cohen Memorial Symposium, Available at: <https://simons.berkeley.edu/talks/tba-4>. 2017 (cit. on p. 1).
- [40] D. Spielman and S. Teng. “Nearly Linear Time Algorithms for Preconditioning and Solving Symmetric, Diagonally Dominant Linear Systems”. In: *SIAM Journal on Matrix Analysis and Applications* 35.3 (2014). Available at <http://arxiv.org/abs/cs/0607105>, pp. 835–885. DOI: 10.1137/090771430 (cit. on pp. 2, 32).
- [41] Shang-Hua Teng. “The Laplacian Paradigm: Emerging Algorithms for Massive Graphs”. In: *Theory and Applications of Models of Computation*. 2010, pp. 2–14 (cit. on p. 2).
- [42] Andreas Tillmann. *Computational aspects of compressed sensing*. Verlag Dr. Hut, 2013 (cit. on p. 1).
- [43] Andreas M Tillmann. “Equivalence of Linear Programming and Basis Pursuit”. In: *PAMM* 15.1 (2015), pp. 735–738 (cit. on p. 1).
- [44] P. M. Vaidya. “Speeding-up linear programming using fast matrix multiplication”. In: *30th Annual Symposium on Foundations of Computer Science*. Oct. 1989, pp. 332–337. DOI: 10.1109/SFCS.1989.63499 (cit. on pp. 2, 3, 4, 5).
- [45] Pravin M. Vaidya. “Solving linear equations with symmetric diagonally dominant matrices by constructing good preconditioners.” Unpublished manuscript UIUC 1990. A talk based on the manuscript was presented at the IMA Workshop on Graph Theory and Sparse Matrix Computation, October 1991, Minneapolis. 1990 (cit. on pp. 2, 21).
- [46] Virginia Vassilevska Williams. “Multiplying Matrices Faster Than Coppersmith-winograd”. In: *STOC*. 2012 (cit. on p. 1).

A Missing Proofs

A.1 Proofs from Section 3

Lemma 3.2. *Function γ_p is as defined above.*

1. For any $p \geq 2$, $t \geq 0$, and $x \in \mathbb{R}$, we have $\gamma_p(t, x) \geq |x|^p$, and $\gamma_p(t, x) \geq \frac{p}{2}t^{p-2}x^2$.
2. It is homogeneous under rescaling of both t and x , i.e., for any $t, \lambda \geq 0, p \geq 1$, and any x we have $\gamma_p(\lambda t, \lambda x) = \lambda^p \gamma_p(t, x)$.
3. For any $t > 0, p \geq 1$ and any x , we have $\gamma_p'(t, x) = p \max\{t, |x|\}^{p-2}x$.

Proof. 1. We have $p \geq 2$. When $|x| \leq t$,

$$\gamma_p(t, x) = \frac{p}{2}t^{p-2}x^2 \geq \frac{p}{2}|x|^p \geq |x|^p.$$

Otherwise,

$$\gamma_p(t, x) = |x|^p + \left(\frac{p}{2} - 1\right) t^p \geq |x|^p.$$

Let $s(x) = \frac{p}{2} t^{p-2} x^2$. At $|x| = t$ we have $\gamma_p(t, t) = s(t)$. Now, $\gamma_p'(t, x) = p|x|^{p-2} x \geq pt^{p-2} x = s'(x)$. This means that for x negative, γ_p decreases faster than s and for x positive, γ_p increases faster than s . The two functions are equal in the range $-t \leq x \leq t$. Therefore, $\gamma_p(t, x) \geq s(x)$ for all x .

2.

$$\begin{aligned} \gamma_p(\lambda t, \lambda x) &= \begin{cases} \frac{p}{2} (\lambda t)^{p-2} (\lambda x)^2 & \text{if } |x| \leq t, \\ |\lambda x|^p + \left(\frac{p}{2} - 1\right) (\lambda t)^p & \text{otherwise.} \end{cases} \\ &= \begin{cases} \lambda^p \frac{p}{2} t^{p-2} x^2 & \text{if } |x| \leq t, \\ \lambda^p |x|^p + \lambda^p \left(\frac{p}{2} - 1\right) t^p & \text{otherwise.} \end{cases} \\ &= \lambda^p \gamma_p(t, x) \end{aligned}$$

3. Taking the derivative of $\gamma_p(t, x)$ with respect to x gives,

$$\frac{d}{dx} \gamma_p = \begin{cases} pt^{p-2} x & \text{if } |x| \leq t, \\ p|x|^{p-1} \cdot \text{sign}(x) & \text{otherwise.} \end{cases}$$

The statement clearly follows. □

Lemma 3.3. For any $p > 1, \Delta \in \mathbb{R}$ and $\lambda \geq 0$, we have,

$$\min\{2, p\} \leq x \frac{\gamma_p'(t, x)}{\gamma_p(t, x)} \leq \max\{2, p\}.$$

This implies,

$$\min\{\lambda^2, \lambda^p\} \gamma_p(t, \Delta) \leq \gamma_p(t, \lambda \Delta) \leq \max\{\lambda^2, \lambda^p\} \gamma_p(t, \Delta).$$

Proof.

$$x \frac{\gamma_p'(t, x)}{\gamma_p(t, x)} = \begin{cases} 2 & \text{if } |x| \leq t, \\ \frac{p|x|^p}{|x|^p + \left(\frac{p}{2} - 1\right) t^p} & \text{otherwise.} \end{cases}$$

Now, when $|x| \geq t$, we have the following. When $p \leq 2$,

$$p = \frac{p|x|^p}{|x|^p} \leq \frac{p|x|^p}{|x|^p + \left(\frac{p}{2} - 1\right) t^p} \leq \frac{p|x|^p}{\frac{p}{2}|x|^p} = 2$$

and when $p \geq 2$,

$$2 = \frac{p|x|^p}{\frac{p}{2}|x|^p} \leq \frac{p|x|^p}{|x|^p + \left(\frac{p}{2} - 1\right) t^p} \leq \frac{p|x|^p}{|x|^p} = \frac{p}{x}$$

The above computations imply that,

$$\min\{2, p\} \leq x \frac{\gamma_p'(t, x)}{\gamma_p(t, x)} \leq \max\{2, p\}.$$

Let $\lambda \geq 1$ and $x \geq 0$. Integrating both sides of the right inequality gives,

$$\begin{aligned} \int_x^{\lambda x} \frac{\gamma_p'(t, x)}{\gamma_p(t, x)} dx &\leq \int_x^{\lambda x} \frac{\max\{2, p\}}{x} dx \\ \Leftrightarrow \log \left(\frac{\gamma_p(t, \lambda x)}{\gamma_p(t, x)} \right) &\leq \log (\lambda)^{\max\{2, p\}} \\ \Leftrightarrow \gamma_p(t, \lambda x) &\leq \lambda^{\max\{2, p\}} \gamma_p(t, x). \end{aligned}$$

Integrating both sides of the left inequality from x to λx gives the required left inequality. Now, let $\lambda \leq 1$. Integrating both sides of the left inequality gives,

$$\begin{aligned} \int_{\lambda x}^x \frac{\gamma_p'(t, x)}{\gamma_p(t, x)} dx &\geq \int_{\lambda x}^x \frac{\min\{2, p\}}{x} dx \\ \Leftrightarrow \log \left(\frac{\gamma_p(t, x)}{\gamma_p(t, \lambda x)} \right) &\geq \log \left(\frac{1}{\lambda} \right)^{\min\{2, p\}} \\ \Leftrightarrow \gamma_p(t, \lambda x) &\leq \lambda^{\min\{2, p\}} \gamma_p(t, x). \end{aligned}$$

Similar to the previous case, integrating both sides of the right inequality from λx to x gives the required left inequality. When $x \leq 0$, the direction of the inequality changes but it gets reversed again after putting limits, since we integrate from λx to x when $\lambda \geq 1$ and x to λx when $\lambda \leq 1$. We thus have,

$$\min\{\lambda^2, \lambda^p\} \gamma_p(t, \Delta) \leq \gamma_p(t, \lambda x) \leq \max\{\lambda^2, \lambda^p\} \gamma_p(t, x)$$

□

Lemma 3.4. *For any $p \geq 2, t \geq 0$ and any x, Δ , we have*

$$\gamma_p(t, x + \Delta) \leq \gamma_p(t, x) + \left| \gamma_p'(t, x) \Delta \right| + p^2 2^{p-3} \max\{t, |x|, |\Delta|\}^{p-2} \Delta^2.$$

Proof. Since $\gamma_p(t, x) = \gamma_p(t, |x|)$, and $\gamma_p(t, x)$ is increasing in x , it suffices to prove the claim for $x, \Delta \geq 0$. We have,

$$\begin{aligned} \gamma_p'(t, x + z) - \gamma_p'(t, x) &= p \max\{t, x + z\}^{p-2} (x + z) - p \max\{t, |x|\}^{p-2} x \\ &= p \max\{t^{p-2} (x + z) - \max\{t, |x|\}^{p-2} x, \\ &\quad (x + z)^{p-1} - \max\{t, |x|\}^{p-2} x\} \\ &\leq p \max\{t^{p-2} (x + z) - t^{p-2} x, (x + z)^{p-1} - x^{p-1}\} \quad (\text{Since } p \geq 2) \\ &\leq p \max\{t^{p-2} z, (p - 1)(x + z)^{p-2} z\} \quad (\text{Using Rolle's theorem}) \\ &\leq p \max\{t^{p-2} z, p(2x)^{p-2} z, p(2z)^{p-2} z\} \\ &\leq p^2 2^{p-2} \max\{t, x, \Delta\}^{p-2} z \quad (\text{Since } z \leq \Delta) \end{aligned}$$

Integrating over $z \in [0, \Delta]$, we get,

$$\gamma_p(x + \Delta) - \gamma_p(x) - \Delta \gamma_p'(x) \leq p^2 2^{p-3} \max\{t, |x|, |\Delta|\}^{p-2} \Delta^2.$$

□

A.2 Proofs from Section 4

Lemma 4.5. *Let $p \in (1, \infty)$. Then for any x and any Δ ,*

$$|x|^p + g\Delta + \frac{p-1}{p2^p}\gamma_p(|x|, \Delta) \leq |x + \Delta|^p \leq |x|^p + g\Delta + 2^p\gamma_p(|x|, \Delta),$$

where $g = p|x|^{p-2}x$ is the derivative of the function $|x|^p$.

Proof. We first show the following two lemmas.

Lemma A.1. *For $|\alpha| \leq 1$ and $p \geq 1$,*

$$1 + \alpha p + \frac{(p-1)}{4}\alpha^2 \leq (1 + \alpha)^p \leq 1 + \alpha p + p2^{p-1}\alpha^2.$$

Proof. Let us first show the left inequality, i.e. $1 + \alpha p + \frac{p-1}{4}\alpha^2 \leq (1 + \alpha)^p$. Define the following function,

$$h(\alpha) = (1 + \alpha)^p - 1 - \alpha p - \frac{p-1}{4}\alpha^2.$$

When $\alpha = 1, -1$, $h(\alpha) \geq 0$. The derivative of h with respect to α is, $h'(\alpha) = p(1 + \alpha)^{p-1} - p - \frac{(p-1)}{2}\alpha$. When $p \geq 2$ and $-1 < \alpha < 1$,

$$\begin{aligned} & \left((1 + \alpha)^{p-2} - 1 \right) \text{sign}(\alpha) \geq 0 \\ & \Rightarrow \left((1 + \alpha)^{p-1} - (1 + \alpha) \right) \text{sign}(\alpha) \geq 0 \\ & \Rightarrow \left(p(1 + \alpha)^{p-1} - p - p\alpha \right) \text{sign}(\alpha) \geq 0 \\ & \Rightarrow \left(p(1 + \alpha)^{p-1} - p - \frac{(p-1)}{2}\alpha \right) \text{sign}(\alpha) \geq 0 \end{aligned}$$

For the last inequality, note that when the product is positive, either both terms are positive or both terms are negative. When both terms are positive, subtracting $(p-1)/2$ instead of p gives a larger positive quantity. When both terms are negative then subtracting $(p-1)/2$ instead of p gives only a smaller quantity, so the inequality holds. This shows that $h'(\alpha)\text{sign}(\alpha) \geq 0$, which means minimum of h is at $h(0) = 0$. Next let us see what happens when $p \leq 2$ and $|\alpha| < 1$.

$$h''(\alpha) = p(p-1)(1 + \alpha)^{p-2} - \frac{p-1}{2} = (p-1) \left(\frac{p}{(1 + \alpha)^{2-p}} - \frac{1}{2} \right) \geq 0$$

This implies that $h'(\alpha)$ is an increasing function of α and α_0 for which $h'(\alpha_0) = 0$ is where h attains its minimum value. The only point where h' is 0 is $\alpha_0 = 0$. This implies $h(\alpha) \geq h(0) = 0$. This concludes the proof of the left inequality. For the right inequality, define:

$$s(\alpha) = 1 + \alpha p + p2^{p-1}\alpha^2 - (1 + \alpha)^p.$$

Note that $s(0) = 0$ and $s(1), s(-1) \geq 0$. We have,

$$s'(\alpha) = p + p2^p\alpha - p(1 + \alpha)^{p-1}$$

Using the mean value theorem for $p \geq 2$ and $\alpha < 0$,

$$\begin{aligned} (1 + \alpha)^{p-1} - 1 &= (p-1)\alpha \cdot (1+z)^{p-2}, z \in (\alpha, 0) \\ &\geq \alpha 2^p. \end{aligned}$$

This implies that $s'(\alpha) \leq 0$ for negative alpha. When $1 > \alpha > 0$, using the convexity of $f(x) = (1+x)^{p-1}$ for $p > 2$, we get,

$$f(\alpha \cdot 1 + (1-\alpha) \cdot 0) \leq \alpha f(1) + (1-\alpha)f(0)$$

which gives us

$$(1 + \alpha)^{p-1} \leq \alpha 2^{p-1} + 1.$$

This implies, $s'(\alpha) \geq 0$ for positive α . The function s is thus increasing for positive α and decreasing for negative α , so it attains the minimum at 0 which is $s(0) = 0$ giving us $s(\alpha) \geq 0$. We now look at the case $p \leq 2$. We have

$$(1 + \alpha)^{p-1} \text{sign}(\alpha) \leq (1 + \alpha) \text{sign}(\alpha).$$

Using this, we get, $s'(\alpha) \text{sign}(\alpha) \geq p|\alpha|(2^p - 1) \geq 0$ which says $s'(\alpha)$ is positive for α positive and negative for α negative. Thus the minima of s is at 0 which is 0. So $s(\alpha) \geq 0$ in this range too. \square

Lemma A.2. For $\beta \geq 1$ and $p \geq 1$, $(\beta - 1)^{p-1} + 1 \geq \frac{1}{2^p} \beta^{p-1}$.

Proof. $(\beta - 1) \geq \frac{\beta}{2}$ for $\beta \geq 2$. So the claim clearly holds for $\beta \geq 2$ since $(\beta - 1)^{p-1} \geq \left(\frac{\beta}{2}\right)^{p-1}$.

When $1 \leq \beta \leq 2$, $1 \geq \frac{\beta}{2}$, so the claim holds since, $1 \geq \left(\frac{\beta}{2}\right)^{p-1}$ \square

We now prove the theorem.

Let $\Delta = \alpha x$. The term $g\Delta = p|x|^{p-1} \text{sign}(x) \cdot \alpha x = \alpha p|x|^{p-1}|x| = \alpha p|x|^p$. Let us first look at the case when $|\alpha| \leq 1$. We want to show,

$$\begin{aligned} |x|^p + \alpha p|x|^p + c\frac{p}{2}|x|^{p-2}|\alpha x|^2 &\leq |x + \alpha x|^p \leq |x|^p + \alpha p|x|^p + C\frac{p}{2}|x|^{p-2}|\alpha x|^2 \\ \Leftrightarrow (1 + \alpha p) + c\frac{p}{2}\alpha^2 &\leq (1 + \alpha)^p \leq (1 + \alpha p) + C\frac{p}{2}\alpha^2. \end{aligned}$$

This follows from Lemma A.1 and the facts $\frac{cp}{2} \leq \frac{p-1}{4}$ and $\frac{Cp}{2} \geq p2^{p-1}$. We next look at the case when $|\alpha| \geq 1$. Now, $\gamma_{|f|}^p(\Delta) = |\Delta|^p + (\frac{p}{2} - 1)|f|^p$. We need to show

$$|x|^p(1 + \alpha p) + \frac{|x|^p(p-1)}{p2^p}(|\alpha|^p + \frac{p}{2} - 1) \leq |x|^p|1 + \alpha|^p \leq |x|^p(1 + \alpha p) + 2^p|x|^p(|\alpha|^p + \frac{p}{2} - 1).$$

When $|x| = 0$ it is trivially true. When $|x| \neq 0$, let

$$h(\alpha) = |1 + \alpha|^p - (1 + \alpha p) - \frac{(p-1)}{p2^p}(|\alpha|^p + \frac{p}{2} - 1).$$

Now, taking the derivative with respect to α we get,

$$h'(\alpha) = p \left(|1 + \alpha|^{p-1} \text{sign}(\alpha) - 1 - \frac{(p-1)}{p2^p} |\alpha|^{p-1} \text{sign}(\alpha) \right).$$

When $\alpha \geq 1$ and $p \geq 2$,

$$h''(\alpha) \geq p(p-1)(1+\alpha)^{p-2} - p \frac{(p-1)}{2^p} \alpha^{p-2} \geq 0.$$

So we have $h'(\alpha) \geq h'(1) \geq 0$. When $p < 2$, we use the mean value theorem to get,

$$\begin{aligned} (1+\alpha)^{p-1} - 1 &= (p-1)\alpha(1+z)^{p-2}, z \in (0, \alpha) \\ &\geq (p-1)\alpha(2\alpha)^{p-2} \\ &\geq \frac{p-1}{2} \alpha^{p-1} \end{aligned}$$

which implies $h'(\alpha) \geq 0$ in this range as well. When $\alpha \leq -1$ it follows from Lemma A.2 that $h'(\alpha) \leq 0$. So the function h is increasing for $\alpha \geq 1$ and decreasing for $\alpha \leq -1$. The minimum value of h is $\min\{h(1), h(-1)\} \geq 0$. It follows that $h(\alpha) \geq 0$ which gives us the left inequality. The other side requires proving,

$$|1 + \alpha|^p \leq 1 + \alpha p + 2^p \left(|\alpha|^p + \frac{p}{2} - 1 \right).$$

Define:

$$s(\alpha) = 1 + \alpha p + 2^p \left(|\alpha|^p + \frac{p}{2} - 1 \right) - |1 + \alpha|^p.$$

The derivative $s'(\alpha) = p + (p2^p|\alpha|^{p-1} - p|1 + \alpha|^{p-1}) \text{sign}(\alpha)$ is non negative for $\alpha \geq 1$ and non positive for $\alpha \leq -1$. The minimum value taken by s is $\min\{s(1), s(-1)\}$ which is non negative. This gives us the right inequality. □

Lemma 4.8. For $\mathbf{x}^{(0)}$ as defined in Definition 4.7, $\|\mathbf{x}^{(0)}\|_p^p \leq m^{(p-2)/2} \text{OPT}$.

Proof. Let \mathbf{x}^* give the OPT. We know that, for any \mathbf{x} ,

$$\|\mathbf{x}\|_p \leq \|\mathbf{x}\|_2 \leq m^{(p-2)/2p} \|\mathbf{x}\|_p.$$

This along with the fact $\|\mathbf{x}^{(0)}\|_2 \leq \|\mathbf{x}^*\|_2$ gives us,

$$\|\mathbf{x}^{(0)}\|_p^p \leq \|\mathbf{x}^{(0)}\|_2^p \leq m^{(p-2)/2} \|\mathbf{x}^*\|_p^p = m^{(p-2)/2} \text{OPT}.$$

□

A.3 Proofs from Section 5

Lemma 5.3. Let $p \in (1, \infty)$ and assume that our current solution \mathbf{x} is not an ε -approximate solution. Let λ be such that $\lambda^{\min\{1, p-1\}} = \frac{p-1}{p4^p}$. For some

$$i \in \left[\log \left(\frac{\varepsilon \|\mathbf{x}^{(0)}\|_p^p}{m^{(p-2)/2}} \right), \log \left(\frac{\|\mathbf{x}^{(0)}\|_p^p}{\lambda} \right) \right],$$

$\alpha(\Delta^*) \in [2^{i-1}, 2^i)$ where Δ^* is the solution of (1).

Proof. Let \mathbf{x}^* denote the optimum solution of (*) and $\mathbf{x}^{(0)}$ be as defined in Definition 4.7. We know that for any \mathbf{x} ,

$$\|\mathbf{x}\|_p \leq \|\mathbf{x}\|_2 \leq m^{(p-2)/2p} \|\mathbf{x}\|_p.$$

This along with the fact $\|\mathbf{x}^{(0)}\|_2 \leq \|\mathbf{x}^*\|_2$ gives us,

$$\|\mathbf{x}^{(0)}\|_p^p \leq \|\mathbf{x}^{(0)}\|_2^p \leq m^{(p-2)/2} \|\mathbf{x}^*\|_p^p.$$

Now from Lemma 4.6 we have,

$$\alpha(\Delta) \leq \frac{1}{\lambda} (\|\mathbf{x}^{(0)}\|_p^p - \|\mathbf{x}^*\|_p^p) \leq \frac{\|\mathbf{x}^{(0)}\|_p^p}{\lambda} (1 - m^{-(p-2)/2}).$$

Let us assume $\alpha(\Delta) \geq \varepsilon \|\mathbf{x}^*\|_p^p \geq \varepsilon m^{-(p-2)/2} \|\mathbf{x}^{(0)}\|_p^p$. If this is not true we already have an ε approximate solution to our problem. We thus have the following bound on α ,

$$\varepsilon m^{-(p-2)/2} \|\mathbf{x}^{(0)}\|_p^p \leq \alpha(\Delta) \leq \frac{\|\mathbf{x}^{(0)}\|_p^p}{\lambda}.$$

This gives us that,

$$2^{\log\left(\frac{\varepsilon \|\mathbf{x}^{(0)}\|_p^p}{m^{(p-2)/2}}\right)} \leq \alpha(\Delta) \leq 2^{\log\left(\frac{\|\mathbf{x}^{(0)}\|_p^p}{\lambda}\right)}.$$

When $p \leq 2$, following a similar proof and using,

$$\|\mathbf{x}\|_2 \leq \|\mathbf{x}\|_p \leq m^{(2-p)/2p} \|\mathbf{x}\|_2,$$

we get,

$$2^{\log\left(\frac{\varepsilon \|\mathbf{x}^{(0)}\|_p^p}{m^{(2-p)/2}}\right)} \leq \alpha(\Delta) \leq 2^{\log\left(\frac{\|\mathbf{x}^{(0)}\|_p^p}{\lambda}\right)},$$

thus concluding the proof of the lemma. thus concluding the proof of the lemma. \square

Lemma 5.4. *Let $p \in (1, \infty)$. Suppose $\alpha(\Delta^*) \in [2^{i-1}, 2^i)$ for some i where Δ^* is the solution of (1). The following program is feasible:*

$$\begin{aligned} \gamma_p(\mathbf{t}, \Delta) &\leq \frac{p}{p-1} 2^{i+p}, \\ \mathbf{g}^T \Delta &= 2^{i-1}, \\ \mathbf{A} \Delta &= 0. \end{aligned} \tag{3}$$

If $\Delta(i)$ is a β -approximate solution to program (3) for this choice of i , then, we can pick $\mu \leq 1$ such that the vector $\mu \Delta(i)$ is an $\Omega_p\left(\beta^{\frac{1}{\min\{p, 2\}-1}}\right)$ -approximate solution to (1).

Proof. Assume that the optimum solution to (1), Δ^* satisfies

$$\alpha(\Delta^*) = \mathbf{g}^T \Delta^* - \frac{p-1}{p2^p} \gamma_p(\mathbf{t}, \Delta^*) \in [2^{i-1}, 2^i),$$

in addition to $\mathbf{A}\Delta^* = 0$. Note that we know that the objective is strictly positive (as 0 is a feasible solution). Since $\gamma_p \geq 0$, we must have,

$$\mathbf{g}^T \Delta^* \geq 2^{i-1},$$

Consider scaling Δ^* by a factor $\lambda > 0$. Since Δ^* is optimal, we must have

$$\left. \frac{d}{d\lambda} \left(\lambda \mathbf{g}^T \Delta^* - \frac{p-1}{p2^p} \gamma_p(\mathbf{t}, \lambda \Delta^*) \right) \right|_{\lambda=1} = 0.$$

Now, from Lemma 3.3, we know that

$$\frac{\min\{2, p\}}{\lambda} \gamma_p(\mathbf{t}, \lambda \Delta^*) \leq \frac{d}{d\lambda} \gamma_p(\mathbf{t}, \lambda \Delta^*) \leq \frac{\max\{2, p\}}{\lambda} \gamma_p(\mathbf{t}, \lambda \Delta^*).$$

Thus, we get,

$$\frac{p-1}{p2^p} \min\{2, p\} \gamma_p(\mathbf{t}, \Delta^*) \leq \mathbf{g}^T \Delta^* \leq \frac{p-1}{p2^p} \max\{2, p\} \gamma_p(\mathbf{t}, \Delta^*),$$

$$\frac{p-1}{p2^p} (\min\{2, p\} - 1) \gamma_p(\mathbf{t}, \Delta^*) \leq \mathbf{g}^T \Delta^* - \frac{p-1}{p2^p} \gamma_p(\mathbf{t}, \Delta^*) \leq \frac{p-1}{p2^p} (\max\{2, p\} - 1) \gamma_p(\mathbf{t}, \Delta^*).$$

Thus, $\gamma_p(\mathbf{t}, \Delta^*) \leq \frac{1}{\frac{p-1}{p2^p} (\min\{2, p\} - 1)} 2^i$, and hence $\mathbf{g}^T \Delta^* \leq \frac{\max\{2, p\}}{\min\{2, p\} - 1} 2^i = \max\{p, 2/p-1\} 2^i$.

Now consider the vector $\Delta = \lambda \Delta^*$, where $\lambda = \frac{2^{i-1}}{\mathbf{g}^T \Delta^*}$. Note that $\lambda \in \left[\min\{1/2p, (p-1)/4\}, 1 \right]$. We have

$$\mathbf{g}^T \Delta = \mathbf{g}^T \left(\frac{2^{i-1}}{\mathbf{g}^T \Delta^*} \Delta^* \right) = 2^{i-1}$$

$$\gamma_p(\mathbf{t}, \Delta) \leq \gamma_p(\mathbf{t}, \lambda \Delta^*) \leq \max\{\lambda^2, \lambda^p\} \gamma_p(\mathbf{t}, \Delta^*) \leq \frac{p2^p}{p-1} 2^i.$$

Thus, Δ is a feasible solution to Program (3). A β -approximate solution $\Delta(i)$ must be such that,

$$\mathbf{A}\Delta(i) = 0,$$

$$\mathbf{g}^T \Delta(i) = 2^{i-1},$$

$$\gamma_p(\mathbf{t}, \Delta(i)) \leq \beta \frac{p2^p}{p-1} 2^i.$$

Now, we consider $\Delta = \mu \Delta(i)$ for some $\mu \leq 1$. We have, $\mathbf{A}\Delta = 0$, and,

$$\begin{aligned} \mathbf{g}^T \Delta - \frac{p-1}{p2^p} \gamma_p(\mathbf{t}, \Delta) &= \mu \mathbf{g}^T \Delta(i) - \frac{p-1}{p2^p} \gamma_p(\mathbf{t}, \mu \Delta(i)) \\ &\geq \mu 2^{i-1} - \max\{\mu^2, \mu^p\} \beta 2^i \end{aligned} \quad \text{(Using Lemma 3.3)}$$

We can pick,

$$\mu \leq \begin{cases} \left(\frac{1}{2\beta p} \right)^{\frac{1}{p-1}} & \text{if } p \leq 2 \\ \frac{1}{4\beta} & \text{if } p \geq 2. \end{cases}$$

In either case, we get,

$$\mathbf{g}^\top \Delta - \frac{p-1}{p2^p} \gamma_p(\mathbf{t}, \Delta) \geq \mu \left(1 - \frac{1}{\min\{2, p\}} \right) 2^{i-1}$$

Since we assumed that the optimum of Program (1) is at most 2^i , this implies that $\mu\Delta(i)$ achieves an objective value for Program (1) that is within an $\Omega_p \left(\beta^{-\frac{1}{\min\{p, 2\}-1}} \right)$ fraction of the optimal. \square

Lemma 5.5. *Let $p \in (1, \infty)$. Let i be such that (3) is feasible. Let*

$$\hat{\mathbf{t}}_j = \begin{cases} m^{-1/p} & \left(\frac{p-1}{p} \right)^{1/p} 2^{-i/p-1} \mathbf{t}_j \leq m^{-1/p}, \\ 1 & \left(\frac{p-1}{p} \right)^{1/p} 2^{-i/p-1} \mathbf{t}_j \geq 1, \\ \left(\frac{p-1}{p} \right)^{1/p} 2^{-i/p-1} \mathbf{t}_j & \text{otherwise.} \end{cases}$$

Note that $m^{-1/p} \leq \hat{\mathbf{t}}_j \leq 1$. Then program (2) with $\mathbf{t} = \hat{\mathbf{t}}$, and

$$\mathbf{c} = \left(\frac{2}{p} \right)^{1/2} \left(\frac{p-1}{p} \right)^{1/p} 2^{i(1-\frac{1}{p})-2},$$

has $\text{OPT} \leq 1$. Let Δ^* be a κ -approximate solution to (2). Then, $\Delta = \left(\frac{p}{2} \right)^{1/2} \left(\frac{p}{p-1} \right)^{1/p} 2^{1+i/p} \Delta^*$ is a $\Omega_p(\kappa)$ -approximate solution to (3).

Proof. We choose i such that (3) is feasible, i.e., there exists Δ such that,

$$\begin{aligned} \gamma_p(\mathbf{t}, \Delta) &\leq \frac{p2^p}{p-1} 2^i, \\ \mathbf{g}^\top \Delta &= 2^{i-1}, \\ \mathbf{A}\Delta &= 0. \end{aligned}$$

Scaling both \mathbf{t} and Δ to $\tilde{\mathbf{t}} = \left(\frac{p-1}{p} \right)^{1/p} 2^{-1-i/p} \mathbf{t}$ and $\tilde{\Delta} = \left(\frac{p-1}{p} \right)^{1/p} 2^{-1-i/p} \Delta$ gives us the following.

$$\begin{aligned} \gamma_p(\tilde{\mathbf{t}}, \tilde{\Delta}) &\leq 1, \\ \mathbf{g}^\top \tilde{\Delta} &= \left(\frac{p-1}{p} \right)^{1/p} 2^{i(1-1/p)-2}, \\ \mathbf{A}\tilde{\Delta} &= 0. \end{aligned}$$

Now, let $\mathbf{t}' = \max\{m^{-1/p}, \tilde{\mathbf{t}}_e\}$. We claim that when $p \geq 2$, $\gamma_p(\mathbf{t}', \tilde{\Delta}) - \gamma_p(\tilde{\mathbf{t}}, \tilde{\Delta}) \leq \frac{p}{2} - 1$. To see this, for a single j , let us look at the difference $\gamma_p(\mathbf{t}'_j, \tilde{\Delta}_j) - \gamma_p(\tilde{\mathbf{t}}_j, \tilde{\Delta}_j)$. If $\tilde{\mathbf{t}}_j \geq m^{-1/p}$ the difference is 0. Otherwise from the proof of Lemma 5 of [5],

$$\gamma_p(\mathbf{t}'_j, \tilde{\Delta}_j) - \gamma_p(\tilde{\mathbf{t}}_j, \tilde{\Delta}_j) \leq \gamma_p(\mathbf{t}'_j, \tilde{\Delta}_j) - |\tilde{\Delta}_j|^p \leq \left(\frac{p}{2} - 1 \right) \left(m^{-1/p} \right)^p.$$

When $p \leq 2$, we claim that $\gamma_p(\tilde{\mathbf{t}}, \tilde{\Delta}) - \gamma_p(\mathbf{t}', \tilde{\Delta}) \leq 1 - \frac{p}{2}$. Again if $\tilde{\mathbf{t}}_j \geq m^{-1/p}$ the difference is 0. Otherwise,

$$\gamma_p(\tilde{\mathbf{t}}_j, \tilde{\Delta}_j) - \gamma_p(\mathbf{t}'_j, \tilde{\Delta}_j) \leq |\tilde{\Delta}_j|^p - \gamma_p(\mathbf{t}'_j, \tilde{\Delta}_j) \leq \left(1 - \frac{p}{2}\right) \left(m^{-1/p}\right)^p.$$

To see the last inequality, when $|\Delta_j| \leq \mathbf{t}'_j$, we require, $|\Delta_j|^p - \frac{p}{2} \mathbf{t}'_j{}^{p-2} \Delta_j^2 \leq \left(1 - \frac{p}{2}\right) \mathbf{t}'_j{}^p$ which is true. When $|\Delta_j| \geq \mathbf{t}'_j$, it directly follows. Summing over all j gives us our claims. We know that $\gamma_p(\tilde{\mathbf{t}}, \tilde{\Delta}) \leq 1$. Thus, $\gamma_p(\mathbf{t}', \tilde{\Delta}) \leq \frac{p}{2}$. Next we set $\hat{\Delta} = \left(\frac{2}{p}\right)^{1/2} \tilde{\Delta}$. Note that $\max\left\{\left(\frac{2}{p}\right)^2, \left(\frac{2}{p}\right)^p\right\} = \left(\frac{2}{p}\right)^2$ for all p . Lemma 3.3 thus implies,

$$\gamma_p(\mathbf{t}', \hat{\Delta}) \leq \left(\frac{2}{p}\right) \gamma_p(\mathbf{t}', \tilde{\Delta}) \leq 1.$$

Define $\hat{\mathbf{t}}_j = \min\{1, \mathbf{t}'_j\}$. Note that $\gamma_p(\hat{\mathbf{t}}, \hat{\Delta}) = \gamma_p(\mathbf{t}', \hat{\Delta})$ since $\gamma_p(\mathbf{t}', \hat{\Delta}) \leq 1$ and as a result we have $\gamma_p(\hat{\mathbf{t}}, \hat{\Delta}) \leq 1$. Observe that $\hat{\Delta}$ is a feasible solution of (2) thus suggesting that for problem (2) $\text{OPT} \leq 1$. Let Δ^* be a κ -approximate solution to (2), i.e.,

$$\gamma_p(\hat{\mathbf{t}}, \Delta^*) \leq \kappa \cdot \text{OPT} \leq \kappa.$$

When $p \geq 2$, γ_p is an increasing function of \mathbf{t} giving us,

$$\gamma_p(\tilde{\mathbf{t}}, \Delta^*) \leq \gamma_p(\mathbf{t}', \Delta^*) = \gamma_p(\hat{\mathbf{t}}, \Delta^*) \leq \kappa.$$

When $p \leq 2$,

$$\gamma_p(\tilde{\mathbf{t}}, \Delta^*) \leq \gamma_p(\mathbf{t}', \Delta^*) + 1 - \frac{p}{2} \leq \kappa + 1$$

This gives,

$$\gamma_p\left(\mathbf{t}, \left(\frac{p}{p-1}\right)^{1/p} 2^{1+i/p} \Delta^*\right) \leq \frac{p2^p}{p-1} 2^i (\kappa + 1)$$

and Lemma 3.3 then implies,

$$\gamma_p\left(\mathbf{t}, \left(\frac{p}{2}\right)^{1/2} \left(\frac{p}{p-1}\right)^{1/p} 2^{1+i/p} \Delta^*\right) \leq \left(\frac{p}{2}\right)^{p/2} \frac{p2^p}{p-1} 2^i (\kappa + 1).$$

Finally, $\Delta = \left(\frac{p}{2}\right)^{1/2} \left(\frac{p}{p-1}\right)^{1/p} 2^{1+i/p} \Delta^*$ satisfies the constraints of (3) and is a $\Omega_p(\kappa)$ approximate solution. \square

Lemma 5.6. *Let $p \geq 2$. For any set of weights \mathbf{w} on the edges, $\sum_e \mathbf{w}_e^{p-2} (\Delta_e^*)^2 \leq \|\mathbf{w}\|_p^{p-2}$.*

Proof. Using Hölder's inequality, we have,

$$\sum_e \mathbf{w}_e^{p-2} (\Delta_e^*)^2 \leq \left(\sum_e ((\Delta_e^*)^2)^{p/2}\right)^{2/p} \left(\sum_e (|\mathbf{w}_e|^{p-2})^{p/(p-2)}\right)^{(p-2)/p}$$

$$\begin{aligned}
&= \left(\sum_e |\Delta_e^*|^p \right)^{2/p} \left(\sum_e |\mathbf{w}_e|^p \right)^{(p-2)/p} \\
&\leq \left(\sum_e |\mathbf{w}_e|^p \right)^{\frac{(p-2)}{p}}, \text{ since } \sum_e (\Delta_e^*)^p \leq 1.
\end{aligned}$$

□

A.4 Proofs from Section 6

Lemma 6.7. *Consider a flow step from Line 13 of Algorithm 4. We have*

$$\frac{\mathbf{r}_e^{(i+1)} - \mathbf{r}_e^{(i)}}{\mathbf{r}_e^{(i)}} \leq (1 + \alpha |\Delta_e|)^{p-2} - 1$$

where Δ is the ℓ_2 minimizer solution produced by the oracle.

Proof. Recall from the setting of resistances from Line 2 of ORACLE (Algorithm 3) that

$$\mathbf{r}_e^{(i)} = \left(m^{1/p} \mathbf{t}_e \right)^{p-2} + \left(\mathbf{w}_e^{(i)} \right)^{p-2}.$$

By Line 13 of Algorithm 4, we have

$$\mathbf{w}_e^{(i+1)} - \mathbf{w}_e^{(i)} = \alpha |\Delta_e|.$$

Substituting this in gives

$$\frac{\mathbf{r}_e^{(i+1)} - \mathbf{r}_e^{(i)}}{\mathbf{r}_e^{(i)}} = \frac{\left(\mathbf{w}_e^{(i)} + \alpha |\Delta_e| \right)^{p-2} - \left(\mathbf{w}_e^{(i)} \right)^{p-2}}{\left(m^{1/p} \mathbf{t}_e \right)^{p-2} + \left(\mathbf{w}_e^{(i)} \right)^{p-2}}.$$

There are two cases to consider:

1. $\mathbf{w}_e^{(i)} \geq m^{1/p} \mathbf{t}_e$.

$$\frac{\mathbf{r}_e^{(i+1)} - \mathbf{r}_e^{(i)}}{\mathbf{r}_e^{(i)}} \leq \frac{\left(\mathbf{w}_e^{(i)} + \alpha |\Delta_e| \right)^{p-2} - \left(\mathbf{w}_e^{(i)} \right)^{p-2}}{\left(\mathbf{w}_e^{(i)} \right)^{p-2}} \leq \left(1 + \frac{\alpha |\Delta_e|}{\mathbf{w}_e^{(i)}} \right)^{p-2} - 1 \leq (1 + \alpha |\Delta_e|)^{p-2} - 1$$

where the last inequality utilizes $\mathbf{w}_e^{(i)} \geq 1$, which is due to the assumption and $m^{1/p} \mathbf{t}_e \geq 1$.

2. $\mathbf{w}_e^{(i)} \leq m^{1/p} \mathbf{t}_e$, then replacing the denominator with the $(m^{1/p} \mathbf{t}_e)^{p-2}$ term and simplifying gives

$$\frac{\mathbf{r}_e^{(i+1)} - \mathbf{r}_e^{(i)}}{\mathbf{r}_e^{(i)}} \leq \left(\frac{\mathbf{w}_e^{(i)}}{m^{1/p} \mathbf{t}_e} + \frac{\alpha |\Delta_e|}{m^{1/p} \mathbf{t}_e} \right)^{p-2} - \left(\frac{\mathbf{w}_e^{(i)}}{m^{1/p} \mathbf{t}_e} \right)^{p-2}.$$

As the function $(z + \theta)^{p-2} - z^{p-2}$ is monotonically increasing when $\theta, p - 2 \geq 0$, we may replace the $\frac{\mathbf{w}_e^{(i)}}{m^{1/p} \mathbf{t}_e}$ by its upper of 1 (given by the assumption) to get

$$\frac{\mathbf{r}_e^{(i+1)} - \mathbf{r}_e^{(i)}}{\mathbf{r}_e^{(i)}} \leq \left(1 + \frac{\alpha |\Delta_e|}{m^{1/p} \mathbf{t}_e}\right)^{p-2} - 1 \leq (1 + \alpha |\Delta_e|)^{p-2} - 1,$$

where the last inequality follows from $m^{1/p} \mathbf{t}_e \geq 1$.

□

B Controlling Φ

Lemma 5.10. *After i flow steps, and k width-reduction steps, provided*

1. $\alpha^p \tau \leq \alpha m^{\frac{p-1}{p}}$, (controls Φ growth in flow-steps)
2. $k \leq \rho^2 m^{2/p} \beta^{-\frac{2}{p-2}}$, (acceptable number of width-reduction steps)

the potential Φ is bounded as follows:

$$\Phi(i, k) \leq \left(p^2 2^p \alpha i + m^{1/p}\right)^p \exp\left(O_p(1) \frac{k}{\rho^2 m^{2/p} \beta^{-\frac{2}{p-2}}}\right).$$

Proof. We prove this claim by induction. Initially, $i = k = 0$, and $\Phi(0, 0) = 0$, and thus, the claim holds trivially. Assume that the claim holds for some $i, k \geq 0$. We will use Φ as an abbreviated notation for $\Phi(i, k)$ below.

Flow Step. For brevity, we let $\gamma_p(\mathbf{w})$ denote $\gamma_p(m^{1/p} \mathbf{t}, \mathbf{w})$, and use \mathbf{w} to denote $\mathbf{w}^{(i, k)}$.

If the next step is a *flow* step,

$$\begin{aligned} \Phi(i+1, k) &= \gamma_p\left(\mathbf{w}^{(i, k)} + \alpha |\Delta|\right) \\ &\leq \gamma_p(\mathbf{w}) + \alpha \left|\Delta^\top \gamma'(\mathbf{w})\right| + p^2 2^{p-3} \alpha^2 \sum_e \max\{m^{1/p} \mathbf{t}, |\mathbf{w}_e|, \alpha \Delta_e\}^{p-2} \Delta_e^2, \quad \text{by Lemma 3.4} \\ &\leq \gamma_p(\mathbf{w}) + p\alpha \gamma_p(\mathbf{w})^{\frac{p-1}{p}} + p\alpha m^{\frac{p-2}{2p}} \gamma_p(\mathbf{w})^{1/2} \\ &\quad + p^2 2^{p-3} \alpha^2 \sum_e \left(\mathbf{r}_e \Delta_e^2 + \alpha^{p-2} \Delta_e^p\right), \quad \text{by Lemma 5.7} \\ &\leq \gamma_p(\mathbf{w}) + p\alpha \gamma_p(\mathbf{w})^{\frac{p-1}{p}} + p\alpha m^{\frac{p-2}{2p}} \gamma_p(\mathbf{w})^{1/2} \\ &\quad + p^2 2^{p-3} \left(\alpha^2 m^{\frac{p-2}{p}} + \alpha^2 \|\mathbf{w}\|_p^{p-2} + \sum_e \alpha^p \Delta_e^p\right), \quad \text{by Lemma 5.7} \\ &\leq \gamma_p(\mathbf{w}) + p\alpha \gamma_p(\mathbf{w})^{\frac{p-1}{p}} + p\alpha m^{\frac{p-2}{2p}} \gamma_p(\mathbf{w})^{1/2} \\ &\quad + p^2 2^{p-3} \left(\alpha^2 m^{\frac{p-2}{p}} + \alpha^2 \|\mathbf{w}\|_p^{p-2} + \alpha m^{\frac{p-1}{p}}\right), \quad \text{by Assumption 1 of this Lemma} \end{aligned}$$

Using $\|\mathbf{w}\|_p \leq \gamma_p(\mathbf{w})$

$$\begin{aligned} &\leq \gamma_p(\mathbf{w}) + p\alpha\gamma_p(\mathbf{w})^{\frac{p-1}{p}} + p\alpha m^{\frac{p-2}{2p}} \gamma_p(\mathbf{w})^{1/2} \\ &\quad + p^2 2^{p-3} \left(\alpha^2 m^{\frac{p-2}{p}} + \alpha^2 \gamma_p(\mathbf{w})^{\frac{p-2}{p}} + \alpha m^{\frac{p-1}{p}} \right) \end{aligned}$$

Recall $\gamma_p(\mathbf{w}) = \Phi(\mathbf{w})$. Letting z denote $\max\{\Phi(\mathbf{w}), m\}^{1/p}$, we have,

$$\begin{aligned} &\leq z^p + p\alpha z^{(p-1)} + p\alpha z^{(p-1)} \\ &\quad + p^2 2^{p-3} \left(\alpha^2 z^{(p-2)} + \alpha^2 z^{p-2} + \alpha z^{(p-1)} \right) \\ &\leq z^p + p^2 2^{p-2} \alpha z^{(p-1)} \\ &\quad + p^2 2^{p-2} \alpha^2 z^{(p-2)} \\ &\leq (z + p^2 2^p \alpha)^p. \end{aligned}$$

From the inductive assumption, we have

$$\begin{aligned} z = \max\{\Phi, m\}^{1/p} &\leq \max \left\{ \left(p^2 2^p \alpha i + m^{1/p} \right)^p \exp \left(O_p(1) \frac{k}{\rho^2 m^{2/p} \beta^{-\frac{2}{p-2}}} \right), m \right\}^{1/p} \\ &= \left(p^2 2^p \alpha k_1 + m^{1/p} \right) \left(\exp \left(O_p(1) \frac{k}{\rho^2 m^{2/p} \beta^{-\frac{2}{p-2}}} \right) \right)^{1/p}. \end{aligned}$$

Thus,

$$\Phi(i+1, k) \leq (z + p^2 2^p \alpha)^p \leq \left(p^2 2^p \alpha (i+1) + m^{1/p} \right)^p \exp \left(O_p(1) \frac{k}{\rho^2 m^{2/p} \beta^{-\frac{2}{p-2}}} \right)$$

proving the inductive claim.

Width Reduction Step. To analyze a width-reduction step, we first observe that, by Lemma 5.7 and the induction hypothesis, which ensures $\|\mathbf{w}^{(i,k)}\|_p^p \leq \Phi \leq O_p(1)m$, and hence $\sum_e \mathbf{r}_e f_e^2 \leq O_p(1)m^{(p-2)/p}$ so we have

$$\sum_{e \in H} \mathbf{r}_e \leq \rho^{-2} \sum_{e \in H} \mathbf{r}_e f_e^2 \leq \rho^{-2} \sum_e \mathbf{r}_e f_e^2 \leq \rho^{-2} O_p(1) m^{(p-2)/p}.$$

Thus, when the next step is a width-reduction step, we have,

$$\begin{aligned} \Phi(i, k+1) &\leq \Phi + O_p(1) \sum_{e \in H} \mathbf{r}_e^{\frac{p}{p-2}} \\ &\leq \Phi + O_p(1) \left(\sum_{e \in H} \mathbf{r}_e \right) \left(\max_{e \in H} \mathbf{r}_e \right)^{\frac{p}{p-2}-1} \end{aligned}$$

$$\leq \Phi + O_p(1) \left(\rho^{-2} m^{\frac{p-2}{p}} \right) \beta^{\frac{2}{p-2}}$$

Letting z denote $\max\{\Phi(\mathbf{w}), m\}^{1/p}$, we have,

$$\leq z \left(1 + O_p(1) \left(\rho^{-2} m^{-\frac{2}{p}} \right) \beta^{\frac{2}{p-2}} \right).$$

Thus,

$$\Phi(i, k+1) \leq z \left(1 + O_p(1) \left(\rho^{-2} m^{-\frac{2}{p}} \right) \beta^{\frac{2}{p-2}} \right) \leq \left(p^2 2^p \alpha i + m^{1/p} \right)^p \exp \left(O_p(1) \frac{k+1}{\rho^2 m^{2/p} \beta^{-\frac{2}{p-2}}} \right)$$

proving the inductive claim. □

C Solving L2 problems

Lemma C.1. *Given an algorithm SOLVER for solving $\mathbf{B}^\top \mathbf{R}^{-1} \mathbf{B} \mathbf{x} = \mathbf{d}$, for a $m \times n$ -fixed matrix \mathbf{B} , a fixed positive diagonal matrix $\mathbf{R} > 0$ and an arbitrary vector \mathbf{d} , there is an algorithm ENHANCEDSOLVER that can solve*

$$\begin{aligned} \min_{\mathbf{f}} \quad & \frac{1}{2} \mathbf{f}^\top \mathbf{R} \mathbf{f} \\ \text{s.t.} \quad & \mathbf{B}^\top \mathbf{f} = 0 \\ & \mathbf{g}^\top \mathbf{f} = z \end{aligned} \tag{8}$$

with one call to SOLVER, two multiplications of \mathbf{B} with a vector, and an additional $O(m+n)$ time, if we assume

$$\mathbf{g}^\top \mathbf{R}^{-1} \mathbf{B} \left(\mathbf{B}^\top \mathbf{R}^{-1} \mathbf{B} \right)^{-1} \mathbf{B}^\top \mathbf{R}^{-1} \mathbf{g} < \mathbf{g}^\top \mathbf{R}^{-1} \mathbf{g}.$$

Proof. Introducing the Lagrangian multipliers \mathbf{v}, a respectively for the constraint $\mathbf{B}^\top \mathbf{f} = 0$, and $\mathbf{g}^\top \mathbf{f} = z$, we can write the Lagrangian as

$$\frac{1}{2} \mathbf{f}^\top \mathbf{R} \mathbf{f} - \mathbf{v}^\top \mathbf{B}^\top \mathbf{f} - a \left(\mathbf{g}^\top \mathbf{f} - z \right).$$

Now, optimizing the Lagrangian with respect to an unconstrained \mathbf{f} , allows us to write

$$\mathbf{f} = \mathbf{R}^{-1} (\mathbf{B} \mathbf{v} + a \mathbf{g}).$$

Plugging this back, we can simplify our Lagrangian as

$$-\frac{1}{2} (\mathbf{B} \mathbf{v} + a \mathbf{g})^\top \mathbf{R}^{-1} (\mathbf{B} \mathbf{v} + a \mathbf{g}) + a z.$$

Optimizing with respect to a , gives us,

$$a = \frac{z - \mathbf{g}^\top \mathbf{R}^{-1} \mathbf{B} \mathbf{v}}{\mathbf{g}^\top \mathbf{R}^{-1} \mathbf{g}}.$$

Plugging this back, gives the Lagrangian as

$$-\frac{1}{2} \mathbf{v}^\top \mathbf{B}^\top \mathbf{R}^{-1} \mathbf{B} \mathbf{v} + \frac{\left(z - \mathbf{g}^\top \mathbf{R}^{-1} \mathbf{B} \mathbf{v}\right)^2}{2 \mathbf{g}^\top \mathbf{R}^{-1} \mathbf{g}}.$$

We let $\tilde{\mathbf{g}}$ denote the vector $\mathbf{B}^\top \mathbf{R}^{-1} \mathbf{g}$ and \mathbf{M} denote the matrix $\mathbf{B}^\top \mathbf{R}^{-1} \mathbf{B}$. Thus, the Lagrangian can be written as,

$$-\frac{1}{2} \mathbf{v}^\top \mathbf{M} \mathbf{v} + \frac{\left(z - \tilde{\mathbf{g}}^\top \mathbf{v}\right)^2}{2 \mathbf{g}^\top \mathbf{R}^{-1} \mathbf{g}}.$$

This implies that the optimal \mathbf{v} is given by the equation

$$\left(\mathbf{M} - \frac{1}{\mathbf{g}^\top \mathbf{R}^{-1} \mathbf{g}} \tilde{\mathbf{g}} \tilde{\mathbf{g}}^\top\right) \mathbf{v} = -\frac{z \tilde{\mathbf{g}}}{\mathbf{g}^\top \mathbf{R}^{-1} \mathbf{g}}.$$

From the condition assumed on \mathbf{g} , we have $\tilde{\mathbf{g}} \mathbf{M}^{-1} \tilde{\mathbf{g}} < \mathbf{g}^\top \mathbf{R}^{-1} \mathbf{g}$. Thus, we can solve this system using the Sherman-Morrison formula as follows,

$$\begin{aligned} \mathbf{v} &= -\left(\mathbf{M}^{-1} + \frac{\mathbf{M}^{-1} \tilde{\mathbf{g}} \tilde{\mathbf{g}}^\top \mathbf{M}^{-1}}{\mathbf{g}^\top \mathbf{R}^{-1} \mathbf{g} - \tilde{\mathbf{g}}^\top \mathbf{M}^{-1} \tilde{\mathbf{g}}}\right) \frac{z \tilde{\mathbf{g}}}{\mathbf{g}^\top \mathbf{R}^{-1} \mathbf{g}} \\ &= -\frac{z \mathbf{M}^{-1} \tilde{\mathbf{g}}}{\mathbf{g}^\top \mathbf{R}^{-1} \mathbf{g} - \tilde{\mathbf{g}}^\top \mathbf{M}^{-1} \tilde{\mathbf{g}}} \end{aligned}$$

The algorithm ENHANCEDSOLVER computes $\tilde{\mathbf{g}}$, and then invokes SOLVER to compute $\mathbf{M}^{-1} \tilde{\mathbf{g}}$. This allows us to compute \mathbf{v} in an additional $O(m+n)$ time. Finally, we can compute $\mathbf{f} = \mathbf{R}^{-1}(\mathbf{B} \mathbf{v} + \mathbf{a} \mathbf{g})$ using another multiplication with \mathbf{B} and an additional $O(m+n)$ time. \square

D General ℓ_2 Resistance Monotonicity

Lemma 5.11. *Assuming the program (4) is feasible, let Δ be an be a solution to the optimization problem (4) with weights \mathbf{r} . Suppose we increase the resistance on each entry to get \mathbf{r}' . Then,*

$$\Psi(\mathbf{r}') \geq \exp\left(\frac{\sum_e \min\left\{1, \frac{r'_e - r_e}{r_e}\right\} r_e \Delta_e^2}{2\Psi(\mathbf{r})}\right) \Psi(\mathbf{r}).$$

Proof. Recall

$$\begin{aligned} \Psi(\mathbf{r}) &= \min_{\Delta} \sum_e r_e \Delta_e^2 \\ \text{s.t.} \quad &\mathbf{A}' \Delta = \mathbf{c} \end{aligned}$$

Letting \mathbf{R} denote the diagonal matrix with \mathbf{r} on its diagonal, we can write the above as

$$\Psi(\mathbf{r}) = \min_{\Delta} \Delta^\top \mathbf{R} \Delta \tag{9}$$

$$\text{s.t. } \mathbf{A}'\Delta = \mathbf{c}$$

Using Lagrangian duality, and noting that strong duality holds, we can write this as

$$\begin{aligned} \Psi(\mathbf{r}) &= \min_{\Delta} \max_{\mathbf{y}} \Delta^{\top} \mathbf{R}\Delta + 2\mathbf{y}^{\top}(\mathbf{c} - \mathbf{A}'\Delta) \\ &= \max_{\mathbf{y}} \min_{\Delta} 2\mathbf{y}^{\top} \mathbf{c} + \Delta^{\top} \mathbf{R}\Delta - 2\mathbf{y}^{\top} \mathbf{A}'\Delta \end{aligned}$$

The minimizing Δ can be found by setting the gradient w.r.t. to this variable to zero. This gives $2\mathbf{R}\Delta - 2\mathbf{y}^{\top} \mathbf{A}' = 0$, so that $\Delta = \mathbf{R}^{-1}(\mathbf{A}')^{\top} \mathbf{y}$. Plugging in this choice of Δ , we arrive at the dual program

$$\Psi(\mathbf{r}) = \max_{\mathbf{y}} 2\mathbf{c}^{\top} \mathbf{y} - \mathbf{y}^{\top} \mathbf{A}' \mathbf{R}^{-1} (\mathbf{A}')^{\top} \mathbf{y} \quad (10)$$

Crucially, strong duality also implies that if Δ^* is an optimal solution of the primal program (9), and \mathbf{y}^* is an optimal solution to the dual then

$$\min_{\Delta} 2(\mathbf{y}^*)^{\top} \mathbf{c} + \Delta^{\top} \mathbf{R}\Delta - 2(\mathbf{y}^*)^{\top} \mathbf{A}'\Delta$$

is optimized at $\Delta = \Delta^*$. This in turn implies the gradient w.r.t. Δ at $\Delta = \Delta^*$ is zero, so that $\Delta^* = \mathbf{R}^{-1}(\mathbf{A}')^{\top} \mathbf{y}^*$. Let \mathbf{a}_e be the e th row of \mathbf{A}' . Then the previous equation tells us that $\Delta_e^* = \frac{1}{r_e} \mathbf{a}_e^{\top} \mathbf{y}^*$. This implies that

$$\forall e. r_e (\Delta_e^*)^2 = \frac{1}{r_e} (\mathbf{a}_e^{\top} \mathbf{y}^*)^2. \quad (11)$$

Consider another program, essentially the same as (10), but with additional scalar valued variable $\theta \in \mathbb{R}$ introduced.

$$\Psi(\mathbf{r}) = \max_{\mathbf{z}, \theta} \theta \cdot 2\mathbf{c}^{\top} \mathbf{z} - \theta^2 \cdot \mathbf{z}^{\top} \mathbf{A}' \mathbf{R}^{-1} (\mathbf{A}')^{\top} \mathbf{z} \quad (12)$$

The two programs (10) and (12) have the same value, since for any \mathbf{y} , the assignment $(\mathbf{z}, \theta) = (\mathbf{y}, 1)$ ensures both objectives take the same value, and conversely for any (\mathbf{z}, θ) , the assignment $\mathbf{y} = \theta \mathbf{z}$ ensures both programs take the same value.

We see that $\mathbf{z} = \mathbf{y}^*$, and $\theta = 1$ is an optimal solution to (12). Hence

$$\left[\frac{d}{d\theta} \left(\theta \cdot 2\mathbf{c}^{\top} \mathbf{y}^* - \theta^2 \cdot (\mathbf{y}^*)^{\top} \mathbf{A}' \mathbf{R}^{-1} (\mathbf{A}')^{\top} \mathbf{y}^* \right) \right]_{\theta=1} = 0$$

Consequently, $\mathbf{c}^{\top} \mathbf{y}^* = (\mathbf{y}^*)^{\top} \mathbf{A}' \mathbf{R}^{-1} (\mathbf{A}')^{\top} \mathbf{y}^*$. Hence $\mathbf{c}^{\top} \mathbf{y}^* = \Psi(\mathbf{r})$. Again, by a scaling argument, this implies that

$$\begin{aligned} 2 - \frac{1}{\Psi(\mathbf{r})} &= \max_{\mathbf{y}} 2\mathbf{c}^{\top} \mathbf{y} - \mathbf{y}^{\top} \mathbf{A}' \mathbf{R}^{-1} (\mathbf{A}')^{\top} \mathbf{y} \\ \text{s.t. } &\mathbf{c}^{\top} \mathbf{y} = 1 \end{aligned}$$

So that

$$\frac{1}{\Psi(\mathbf{r})} = \min_{\mathbf{y}} \mathbf{y}^{\top} \mathbf{A}' \mathbf{R}^{-1} (\mathbf{A}')^{\top} \mathbf{y} \quad (13)$$

$$\text{s.t.} \quad \mathbf{c}^\top \mathbf{y} = 1$$

Note that one optimal assignment for the program (13) is $\tilde{\mathbf{y}} = \frac{\mathbf{y}^*}{\Psi(\mathbf{r})}$. Also observe that if we consider the program (13) with \mathbf{r}' instead of \mathbf{r} as the resistances, then $\tilde{\mathbf{y}} = \frac{\mathbf{y}^*}{\Psi(\mathbf{r})}$ is still a feasible solution. Hence, using the observation (11), we get

$$\frac{1}{\Psi(\mathbf{r}')} \leq \frac{1}{\Psi(\mathbf{r})^2} (\mathbf{y}^*)^\top \mathbf{A}' (\mathbf{R}')^{-1} (\mathbf{A}')^\top \mathbf{y}^* = \frac{1}{\Psi(\mathbf{r})^2} \sum_e \frac{r_e}{r'_e} r_e (\Delta_e^*)^2$$

Factoring out the $\Psi(\mathbf{r})$ term gives

$$\frac{1}{\Psi(\mathbf{r}')} \leq \frac{1}{\Psi(\mathbf{r})} \left(1 - \frac{\sum_e \left(1 - \frac{r_e}{r'_e}\right) r_e (\Delta_e^*)^2}{\Psi(\mathbf{r})} \right) \leq \exp \left(- \frac{\sum_e \left(1 - \frac{r_e}{r'_e}\right) r_e (\Delta_e^*)^2}{\Psi(\mathbf{r})} \right).$$

Now consider the term $1 - \frac{r_e}{r'_e}$: if $r'_e \geq 2r_e$, then it is at least $1/2$. Otherwise, it can be rearranged to

$$\frac{r'_e - r_e}{r'_e} \geq \frac{r'_e - r_e}{2r_e}.$$

So in either case, we have

$$\frac{1}{\Psi(\mathbf{r}')} \leq \frac{1}{\Psi(\mathbf{r})} \exp \left(- \frac{\sum_e \min \left\{ 1, \frac{r'_e - r_e}{r_e} \right\} r_e (\Delta_e^*)^2}{2\Psi(\mathbf{r})} \right),$$

which upon rearranging gives the desired result. □