
Iterative Hard Thresholding with Adaptive Regularization: Sparser Solutions Without Sacrificing Runtime

Kyriakos Axiotis¹ Maxim Sviridenko²

Abstract

We propose a simple modification to the iterative hard thresholding (IHT) algorithm, which recovers asymptotically sparser solutions as a function of the condition number. When aiming to minimize a convex function $f(\mathbf{x})$ with condition number κ subject to \mathbf{x} being an s -sparse vector, the standard IHT guarantee is a solution with relaxed sparsity $O(s\kappa^2)$, while our proposed algorithm, *regularized IHT*, returns a solution with sparsity $O(s\kappa)$. Our algorithm significantly improves over ARHT (Axiotis & Sviridenko, 2021b) which also finds a solution of sparsity $O(s\kappa)$, as it does not require re-optimization in each iteration (and so is much faster), is deterministic, and does not require knowledge of the optimal solution value $f(\mathbf{x}^*)$ or the optimal sparsity level s . Our main technical tool is an *adaptive regularization* framework, in which the algorithm progressively learns the weights of an ℓ_2 regularization term that will allow convergence to sparser solutions. We also apply this framework to low rank optimization, where we achieve a similar improvement of the best known condition number dependence from κ^2 to κ .

1. Introduction

Sparse optimization is the task of optimizing a function f over s -sparse vectors, i.e. those with at most s non-zero entries. Examples of such optimization problems arise in machine learning, with the goal to make models smaller for efficiency, generalization, or interpretability reasons, and compressed sensing, where the goal is to recover an s -sparse signal from a small number of measurements. A closely related problem is *low rank optimization*, where the sparsity constraint is instead placed on the spectrum of

¹MIT ²Yahoo! Research. Correspondence to: Kyriakos Axiotis <kaxiotis@mit.edu>, Maxim Sviridenko <sviri@yahooinc.com>.

the solution (which is a matrix). This problem is central in matrix factorization, recommender systems, robust principal components analysis, among other tasks. More generally, *structured sparsity* constraints have the goal of capturing the special structure of a particular task by restricting the set of solutions to those that are “simple” in an appropriate sense. Examples include group sparsity, tree- and graph-structured sparsity. For more on generalized sparsity measures see e.g. (Schmidt, 2018).

Among the huge number of algorithms that have been developed for the sparse optimization problems, three stand out as the most popular ones:

- The **LASSO** (Tibshirani, 1996), which works by relaxing the ℓ_0 (sparsity) constraint to an ℓ_1 constraint, thus convexifying the problem.
- **Orthogonal matching pursuit** (OMP) (Pati et al., 1993), which works by building the solution greedily in an incremental fashion.
- **Iterative hard thresholding** (IHT) (Blumensath & Davies, 2009), which performs projected gradient descent on the set of sparse solutions.

Among these, IHT is generally the most efficient, since it has essentially no overhead over plain gradient descent, making it the tool of choice for large-scale applications.

1.1. Iterative Hard Thresholding (IHT)

Consider the *sparse convex optimization* problem

$$\min_{\|\mathbf{x}\|_0 \leq s} f(\mathbf{x}), \quad (1)$$

where f is convex and $\|\mathbf{x}\|_0$ is the number of non-zero entries in the vector \mathbf{x} , i.e. the sparsity of \mathbf{x} . IHT works by repeatedly performing the following iteration

$$\mathbf{x}^{t+1} = H_{s'}(\mathbf{x}^t - \eta \cdot \nabla f(\mathbf{x}^t)), \quad (2)$$

where $H_{s'}$ is the *hard thresholding operator* that zeroes out all but the top s' entries, for some (potentially relaxed) sparsity level s' , and $\eta > 0$ is the step size.

As (1) is known to be NP-hard (Natarajan, 1995) and even hard to approximate (Foster et al., 2015), an extra assumption needs to be made for the performance of the algorithm to be theoretically evaluated in a meaningful way. The most common assumption is that the (restricted) condition number of f is bounded by κ (or the restricted isometry property constant is bounded by δ (Candes, 2008)), but other assumptions have been studied as well, such as incoherence (Donoho & Elad, 2003) and weak supermodularity (Liberty & Sviridenko, 2017). The performance is then measured in terms of the sparsity s' of the returned solution, as well as its error (value of f).

As it is known (Jain et al., 2014), IHT is guaranteed to return an $s' = O(s\kappa^2)$ -sparse solution \mathbf{x} with $f(\mathbf{x}) \leq f(\mathbf{x}^*) + \varepsilon$. In fact, as we show in Section E, the κ^2 factor cannot be improved in the analysis. Recently, (Axiotis & Sviridenko, 2021b) presented an algorithm called ARHT, which improves the sparsity to $s' = O(s\kappa)$. However, their algorithm is much less efficient than IHT, for many reasons. So the question emerges:

Is there a sparse convex optimization algorithm that returns $O(s\kappa)$ -sparse solutions, but whose runtime efficiency is comparable to IHT?

The main contribution of our work is to show that this goal can be achieved, and done so by a surprisingly simple tweak to IHT.

1.2. Reconciling Sparsity and Efficiency: Regularized IHT

Our main result is the following theorem, which states that running IHT on an *adaptively regularized* objective function returns $O(s\kappa)$ -sparse solutions that are ε -optimal in function value, while having no significant runtime overhead over plain gradient descent.

Theorem 1.1 (Regularized IHT). *Let $f \in \mathbb{R}^n \rightarrow \mathbb{R}$ be a convex function that is β -smooth and α -strongly convex¹, with condition number $\kappa = \beta/\alpha$, and \mathbf{x}^* be an (unknown) s -sparse solution. Then, running Algorithm 1 with $\eta = (2\beta)^{-1}$ and $c = s'/(4T)$ for*

$$T = O\left(\kappa \log \frac{f(\mathbf{x}^0) + (\beta/2) \|\mathbf{x}^0\|_2^2 - f(\mathbf{x}^*)}{\varepsilon}\right)$$

iterations starting from an arbitrary $s' = O(s\kappa)$ -sparse solution \mathbf{x}^0 , the algorithm returns an s' -sparse solution \mathbf{x}^T such that $f(\mathbf{x}^T) \leq f(\mathbf{x}^) + \varepsilon$. Furthermore, each iteration requires $O(1)$ evaluations of f , ∇f , and $O(n)$ additional time.*

¹The theorem also holds if the smoothness and strong convexity constants are replaced by $(s' + s)$ -restricted smoothness and strong convexity constants.

To achieve this result, we significantly refine and generalize the *adaptive regularization* technique of (Axiotis & Sviridenko, 2021b). This refined version fixes many of the shortcomings of the original, by (i) not requiring *re-optimization* in every iteration (a relic of OMP-style algorithms), (ii) taking $\tilde{O}(\kappa)$ instead of $\tilde{O}(s\kappa)$ iterations, (iii) being *deterministic*, (iv) not requiring knowledge of the optimal function value $f(\mathbf{x}^*)$ thus avoiding the overhead of an outer binary search, and (v) being more easily generalizable to other settings, like low rank minimization.

In short, our main idea is to run IHT on a regularized function

$$g(\mathbf{x}) = f(\mathbf{x}) + (\beta/2) \|\mathbf{x}\|_{\mathbf{w},2}^2,$$

where $\|\mathbf{x}\|_{\mathbf{w},2}^2 = \sum_{i=1}^n w_i x_i^2$ and \mathbf{w} are non-negative weights. These weights change dynamically during the algorithm, in a way that depends on the value of \mathbf{x} . The effect is that now the IHT step will instead be given by

$$\mathbf{x}^{t+1} = H_{s'}\left(\left(\mathbf{1} - 0.5\mathbf{w}^t\right) \mathbf{x}^t - \eta \cdot \nabla f(\mathbf{x}^t)\right),$$

which is almost the same as (2), except that it has an extra term that biases the solution towards $\mathbf{0}$. Additionally, in each step the weights \mathbf{w}^t are updated based on the current solution as

$$w_i^{t+1} = \left(w_i^t \cdot \left(\mathbf{1} - c \cdot \frac{w_i^t (x_i^t)^2}{\|\mathbf{x}^t\|_{\mathbf{w}^t,2}^2}\right)\right)_{\geq 1/2}$$

for some parameter $c > 0$, where $(\cdot)_{\geq 1/2}$ denotes zeroing out all the entries that are $< 1/2$ and keeping the others intact.

In Section 3, we will go over the central ideas of our refined adaptive regularization technique, and also explain how it can be extended to deal with more general sparsity measures.

1.3. Beyond Sparsity: Low Rank Optimization

As discussed, our new techniques transfer to the problem of minimizing a convex function under a rank constraint. In particular, we prove the following theorem:

Theorem 1.2 (Adaptive Regularization for Low Rank Optimization). *Let $f \in \mathbb{R}^{m \times n} \rightarrow \mathbb{R}$ be a convex function with condition number κ and consider the low rank minimization problem*

$$\min_{\text{rank}(\mathbf{A}) \leq r} f(\mathbf{A}). \quad (3)$$

For any error parameter $\varepsilon > 0$, there exists a polynomial time algorithm that returns a matrix \mathbf{A} with $\text{rank}(\mathbf{A}) \leq O\left(r \left(\kappa + \log \frac{f(\mathbf{O}) - f(\mathbf{A}^)}{\varepsilon}\right)\right)$ and $f(\mathbf{A}) \leq f(\mathbf{A}^*) + \varepsilon$, where \mathbf{O} is the all-zero matrix and \mathbf{A}^* is any rank- r matrix.*

This result can be compared to the Greedy algorithm of (Axiotis & Sviridenko, 2021a), which works by incrementally adding a rank-1 component to the solution and achieves rank $O(r\kappa \log \frac{f(\mathbf{O})-f(\mathbf{A}^*)}{\varepsilon})$, as well as their Local Search algorithm, which works by simultaneously adding a rank-1 component and removing another, and achieves rank $O(r\kappa^2)$. In contrast, our Theorem 1.2 returns a solution with rank $O\left(r\left(\kappa + \log \frac{f(\mathbf{O})-f(\mathbf{A}^*)}{\varepsilon}\right)\right)$.

1.4. Related Work

The sparse optimization and compressed sensing literature has a wealth of different algorithms and analyses. Examples include the seminal paper of (Candes, 2008) on recovery with LASSO and followup works (Foucart, 2010), the CoSaMP algorithm (Needell & Tropp, 2009), orthogonal matching pursuit and variants (Natarajan, 1995; Shalev-Shwartz et al., 2010; Jain et al., 2011; Axiotis & Sviridenko, 2021b) iterative hard thresholding (Blumensath & Davies, 2009; Jain et al., 2014), hard thresholding pursuit (Foucart, 2011; Yuan et al., 2016; Shen & Li, 2017a;b), partial hard thresholding (Jain et al., 2017), and message passing algorithms (Donoho et al., 2009). For a survey on compressed sensing, see (Boche et al., 2015; Foucart & Rauhut, 2017).

A family of algorithms that is closely related to IHT are Frank-Wolfe (FW) methods (Frank et al., 1956), which have been used for dealing with generalized sparsity constraints (Jaggi, 2013). The basic version can be viewed as a variant of OMP without re-optimization in each iteration. Block-FW methods are more resemblant of IHT without the projection step, see e.g. (Allen-Zhu et al., 2017) for an application to the low rank minimization problem.

(Liu & Foygel Barber, 2020) presented an interesting connection between hard and soft thresholding algorithms by studying a concavity property of the thresholding operator, and proposed new thresholding operators.

Recently it has been shown (Peste et al., 2021) that IHT can be guaranteed to work for sparse optimization of *non-convex* functions, under appropriate assumptions. In particular, (Peste et al., 2021) studies a stochastic version of IHT for sparse deep learning problems, from both a theoretical and practical standpoint.

2. Background

Notation. We denote $[n] = \{1, 2, \dots, n\}$. We will use **bold** to refer to vectors or matrices. We denote by $\mathbf{0}$ the all-zero vector, $\mathbf{1}$ the all-one vector, \mathbf{O} the all-zero matrix, and by \mathbf{I} the identity matrix (with dimensions understood from the context). Additionally, we will denote by $\mathbf{1}_i$ the i -th basis vector, i.e. the vector that is 0 everywhere except at position i .

In order to ease notation and where not ambiguous for two vectors $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$, we denote by $\mathbf{x}\mathbf{y} \in \mathbb{R}^n$ a vector with elements $(\mathbf{x}\mathbf{y})_i = x_i y_i$, i.e. the element-wise multiplication of two vectors \mathbf{x} and \mathbf{y} . In contrast, we denote their inner product by $\langle \mathbf{x}, \mathbf{y} \rangle$ or $\mathbf{x}^\top \mathbf{y}$. Similarly, $\mathbf{x}^2 \in \mathbb{R}^n$ will be the element-wise square of vector \mathbf{x} .

Restrictions and Thresholding. For any vector $\mathbf{x} \in \mathbb{R}^n$ and set $S \subseteq [n]$, we denote by \mathbf{x}_S the vector that results from \mathbf{x} after zeroing out all the entries except those in positions given by indices in S . For any $t \in \mathbb{R}$ and $\mathbf{x} \in \mathbb{R}^n$, we denote by $\mathbf{x}_{\geq t}$ the vector that results from setting all the entries of \mathbf{x} that are less than t to 0. For a function $f \in \mathbb{R}^n \rightarrow \mathbb{R}$, its gradient $\nabla f(\mathbf{x})$, and a set of indices $S \subseteq [n]$, we denote $\nabla_S f(\mathbf{x}) = (\nabla f(\mathbf{x}))_S$. We define the *thresholding operator* $H_s(\mathbf{x})$ for any vector \mathbf{x} as \mathbf{x}_S , where S are the s entries of \mathbf{x} with largest absolute value (breaking ties arbitrarily). We override the thresholding operator $H_r(\mathbf{A})$ when the argument is a matrix \mathbf{A} , defining $H_r(\mathbf{A}) = \mathbf{U} \text{diag}(H_r(\boldsymbol{\lambda})) \mathbf{V}^\top$, where $\mathbf{U} \text{diag}(\boldsymbol{\lambda}) \mathbf{V}^\top$ is the singular value decomposition of \mathbf{A} , i.e. $H_r(\mathbf{A})$ only keeps the top r singular components of \mathbf{A} .

Norms and Inner Products. For any $p \in (0, \infty)$ and weight vector $\mathbf{w} \geq \mathbf{0}$, we define the weighted ℓ_p norm of a vector $\mathbf{x} \in \mathbb{R}^n$ as:

$$\|\mathbf{x}\|_{p, \mathbf{w}} = \left(\sum_i w_i x_i^p \right)^{1/p}.$$

For $p = 0$, we denote $\|\mathbf{x}\|_0 = |\{i \mid x_i \neq 0\}|$ to be the *sparsity* of \mathbf{x} . For $p = \infty$, we denote $\|\mathbf{x}\|_\infty = \max_i |x_i|$ to be the maximum absolute value of \mathbf{x} .

For a matrix $\mathbf{A} \in \mathbb{R}^{m \times n}$, we let $\|\mathbf{A}\|_2$ be its spectral norm, $\|\mathbf{A}\|_F$ be its Frobenius norm, and $\|\mathbf{A}\|_*$ be its nuclear norm (i.e. sum of singular values). For any $\mathbf{B} \in \mathbb{R}^{m \times n}$, we denote the Frobenius inner product as $\langle \mathbf{A}, \mathbf{B} \rangle = \text{Tr}[\mathbf{A}^\top \mathbf{B}]$.

Smoothness, strong convexity, condition number. A differentiable function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is called *convex* if for any $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$ we have $f(\mathbf{y}) \geq f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle$. Furthermore, f is called β -*smooth* for some real number $\beta > 0$ if for any $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$ we have $f(\mathbf{y}) \leq f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle + (\beta/2) \|\mathbf{y} - \mathbf{x}\|_2^2$ and α -*strongly convex* for some real number $\alpha > 0$ if for any $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$ we have $f(\mathbf{y}) \geq f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle + (\alpha/2) \|\mathbf{y} - \mathbf{x}\|_2^2$. We call $\kappa := \beta/\alpha$ the *condition number* of f . If f is only β -smooth along s -sparse directions (i.e. only for $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$ such that $\|\mathbf{y} - \mathbf{x}\|_0 \leq s$), then we call f β -smooth *at sparsity level* s and denote the smallest such β by β_s and call it the *restricted smoothness constant* (at sparsity level s). We analogously define the *restricted strong convexity constant* α_s , as well as the *s-restricted condition number* $\kappa_s := \beta_s/\alpha_s$.

Projections. Given a subspace \mathcal{V} , we will denote the orthogonal projection onto \mathcal{V} as $\mathbf{I}_{\mathcal{V}}$. In particular, for any matrix $\mathbf{A} \in \mathbb{R}^{m \times n}$ we denote by $\text{im}(\mathbf{A}) = \{\mathbf{A}\mathbf{x} \mid \mathbf{x} \in \mathbb{R}^n\}$ the *image* of \mathbf{A} and by $\text{ker}(\mathbf{A}) = \{\mathbf{x} \mid \mathbf{A}^\top \mathbf{x} = \mathbf{0}\}$ the *kernel* of \mathbf{A} . Therefore, $\mathbf{I}_{\text{im}(\mathbf{A})} = \mathbf{A} \left(\mathbf{A}^\top \mathbf{A} \right)^+ \mathbf{A}^\top$ is the orthogonal projection onto the image of \mathbf{A} and $\mathbf{I}_{\text{ker}(\mathbf{A}^\top)} = \mathbf{I} - \mathbf{I}_{\text{im}(\mathbf{A})}$ the orthogonal projection onto the kernel of \mathbf{A}^\top , where $(\cdot)^+$ denotes the matrix pseudoinverse.

3. The Adaptive Regularization Method

Consider the sparse optimization problem

$$\min_{\|\mathbf{x}\|_0 \leq s} f(\mathbf{x}) \quad (4)$$

on a convex function f with condition number at most κ , and an optimal solution \mathbf{x}^* that is supported on the set of indices $S^* \subseteq [n]$.

The main hurdle towards solving this problem is that it is NP hard. Therefore, it is common to relax it by a factor depending on κ . In fact, IHT requires relaxing the sparsity constraint by a factor of $O(\kappa^2)$ (i.e. $\|\mathbf{x}\|_0 \leq O(s\kappa^2)$), in order to return a near-optimal solution. Also, the κ^2 factor is tight for IHT (see Appendix E).

Remark We state all our results in terms of the condition number κ , even though the statements can be strengthened to depend on the *restricted* condition number $\kappa_{s'+s}$, specifically the condition number restricted on $(s' + s)$ -sparse directions. We state our results in this weaker form for clarity of presentation.

3.1. Regularized IHT

Perhaps surprisingly, there is a way to *regularize* the objective by a weighted ℓ_2 norm so that running IHT on the new objective will only require relaxing the sparsity by $O(\kappa)$:

$$\min_{\|\mathbf{x}\|_0 \leq s} f(\mathbf{x}) + (\beta/2) \|\mathbf{x}\|_{w,2}^2. \quad (5)$$

One way to do this is by setting the weights w to be 1 everywhere except in the indices from S^* , where it is set to 0. An inquisitive reader will protest that this is not a very useful statement, since it requires knowledge of S^* , which was our goal to begin with. In fact, we could just as easily have used the regularizer $(\beta/2) \|\mathbf{x} - \mathbf{x}^*\|_2^2$, thus penalizing everything that is far from the optimum!

3.2. Learning Weights

Our main contribution is to show that the optimal weights w can in fact be *learned* in the duration of the algorithm². More precisely, consider running IHT starting from the setting of $w = 1$. The regularized objective (5) is now $O(1)$ -conditioned, which is great news. On the other hand, (5) is not what we set out to minimize. In other words, even though this approach might work great for minimizing (5), it might (and generally will) fail to achieve sufficient decrease in (4)—one could view this as the algorithm getting trapped in a local minimum.

Our main technical tool is to characterize these local minima, by showing that they can only manifest themselves if the current solution \mathbf{x} satisfies the following condition:

$$\|\mathbf{x}_{S^*}\|_{w,2}^2 \geq \Omega(\kappa^{-1}) \|\mathbf{x}\|_{w,2}^2. \quad (6)$$

In words, this means that a significant fraction of the mass of the current solution lies in the support S^* of the optimal solution. Interestingly, this gives us enough information based on which to update the regularization weights w in a way that the sum of weights in S^* drops fast enough compared to the total sum of weights. This implies that the vector w moves in a direction that correlates with the direction of the optimal weight vector.

These are the core ideas needed to bring the sparsity overhead of IHT from $O(\kappa^2)$ down to $O(\kappa)$.

3.3. Beyond Sparsity: Learning Subspaces

One can summarize the approach of the previous section in the following more general way: If we know that the optimal solution \mathbf{x}^* lies in a particular low-dimensional subspace (in our case this was the span of $\mathbf{1}_i$ for all $i \in S^*$), then we can define a regularization term that penalizes all the solutions based on their distance to that subspace. Of course, this subspace is unknown to us, but we can try to adaptively modify the regularization term every time the algorithm gets stuck, just as we did in the previous section.

More concretely, given a collection \mathcal{A} of unit vectors from \mathbb{R}^n (commonly called *atoms*), we define the following problem:

$$\min_{\text{rank}_{\mathcal{A}}(\mathbf{x}) \leq r} f(\mathbf{x}), \quad (7)$$

where $\text{rank}_{\mathcal{A}}(\mathbf{x})$ is the smallest number of vectors from \mathcal{A} such that \mathbf{x} can be written as their linear combination.

²The idea of adaptively learning regularization weights looks on the surface similar to adaptive gradient algorithms such as AdaGrad (Duchi et al., 2011). An important difference is that these algorithms regularize the function around the current solution, while we regularize it around the origin. Still, this is a potentially intriguing connection that deserves to be investigated further.

We can pick $\mathcal{A} = \{\mathbf{1}_1, \mathbf{1}_2, \dots, \mathbf{1}_n\}$ to obtain the sparse optimization problem, $\mathcal{A} = \{\text{vec}(\mathbf{u}\mathbf{v}^\top) \mid \|\mathbf{u}\|_2 = \|\mathbf{v}\|_2 = 1\}$ for the low rank minimization problem, and other choices of \mathcal{A} can capture more sophisticated problem constraints such as graph structure. Defining an IHT variant for these more general settings is usually straightforward, although the analysis for even obtaining a rank overhead of $O(\kappa^2)$ does not trivially follow and depends on the structure of \mathcal{A} .

So, how would a regularizer look in this more general setting? Given our above discussion, it is fairly simple to deduce it. Consider a decomposition of \mathbf{x} as the sum of rank-1 components from $\mathcal{A} = \{\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_{|\mathcal{A}|}\}$:

$$\mathbf{x} = \sum_{i \in S} \mathbf{a}_i,$$

where $\text{rank}_{\mathcal{A}}(\mathbf{a}_i) = 1$, and let $L^* = \text{span}(\{\mathbf{a}_i \mid i \in S^*\})$ be a low-dimensional subspace that contains the optimal solution and L^\perp_* is its complement. We can then define the regularizer

$$\Phi^*(\mathbf{x}) = (\beta/2) \sum_{i \in S} \|\mathbf{I}_{L^\perp_*} \mathbf{a}_i\|_2^2,$$

where $\mathbf{I}_{L^\perp_*}$ is the orthogonal projection onto the subspace perpendicular to L^* —in other words $\|\mathbf{I}_{L^\perp_*} \mathbf{a}_i\|_2$ is the ℓ_2 distance from \mathbf{a}_i to L^* . An equivalent but slightly more concise way is to write:

$$\Phi^*(\mathbf{x}) = (\beta/2) \left\langle \mathbf{I}_{L^\perp_*}, \sum_{i \in S} \mathbf{a}_i \mathbf{a}_i^\top \right\rangle.$$

Then, we can replace the unknown projection matrix $\mathbf{I}_{L^\perp_*}$ by a weight matrix \mathbf{W} initialized at \mathbf{I} , and proceed by adaptively modifying \mathbf{W} as we did in the previous section.

It should be noted that the full analysis of this framework is not automatic for general \mathcal{A} , and there are several technical challenges that arise depending on the choice of \mathcal{A} . In particular, it does not directly apply to the low rank minimization case, and we end up using a different choice of regularizer. However, the discussion in this section should serve as a basic framework for improving the IHT analysis in more general settings, as in particular it did to motivate the low rank optimization analysis that we will present in Section 5.

4. Sparse Optimization Using Regularized IHT

The main result of this section is an efficient algorithm for sparse optimization of convex functions that, even though is a slight modification of IHT, improves the sparsity by an $O(\kappa)$ factor, where κ is the condition number. The regularized IHT algorithm is presented in Algorithm 1 and its

Algorithm 1 Regularized IHT

\mathbf{x}^0 : initial s' -sparse solution
 $\mathbf{w}^0 = \mathbf{1}$: initial regularization weights
 η : step size, T : #iterations
 c : weight step size
for $t = 0 \dots T - 1$ **do**
 $\mathbf{x}^{t+1} = H_{s'}((\mathbf{1} - 0.5\mathbf{w}^t)\mathbf{x}^t - \eta \cdot \nabla f(\mathbf{x}^t))$
 $\mathbf{w}^{t+1} = \left(\mathbf{w}^t - c \cdot (\mathbf{w}^t \mathbf{x}^t)^2 / \|\mathbf{x}^t\|_{\mathbf{w}^t, 2}^2 \right)_{\geq 1/2}$
 if $f(\mathbf{x}^{t+1}) + (4\eta)^{-1} \|\mathbf{x}^{t+1}\|_{\mathbf{w}^{t+1}, 2}^2 > f(\mathbf{x}^t) + (4\eta)^{-1} \|\mathbf{x}^t\|_{\mathbf{w}^{t+1}, 2}^2$ **then**
 $\mathbf{x}^{t+1} = \mathbf{x}^t$ {In practice there is no need to perform this step.}
 end if
end for

analysis is in Theorem 1.1, whose proof can be found in Appendix B.

Theorem 1.1 (Regularized IHT). Let $f \in \mathbb{R}^n \rightarrow \mathbb{R}$ be a convex function that is β -smooth and α -strongly convex, with condition number $\kappa = \beta/\alpha$, and \mathbf{x}^* be an (unknown) s -sparse solution. Then, running Algorithm 1 with $\eta = (2\beta)^{-1}$ and $c = s'/(4T)$ for

$$T = O\left(\kappa \log \frac{f(\mathbf{x}^0) + (\beta/2) \|\mathbf{x}^0\|_2^2 - f(\mathbf{x}^*)}{\varepsilon}\right)$$

iterations starting from an arbitrary $s' = O(s\kappa)$ -sparse solution \mathbf{x}^0 , the algorithm returns an s' -sparse solution \mathbf{x}^T such that $f(\mathbf{x}^T) \leq f(\mathbf{x}^*) + \varepsilon$. Furthermore, each iteration requires $O(1)$ evaluations of f , ∇f , and $O(n)$ additional time.

The main ingredient for proving Theorem 1.1 is Lemma 4.1, which states that each step of the algorithm either makes substantial (multiplicative) progress in an appropriately regularized function $f(\mathbf{x}) + (\beta/2) \|\mathbf{x}\|_{\mathbf{w}, 2}^2$, or a significant fraction of the mass of \mathbf{x}^2 lies in S^* , which is the support of the target solution. This latter condition allows us to adapt the weights \mathbf{w} in order to obtain a new regularization function that penalizes the target solution less. The proof of the lemma can be found in Appendix C.

Lemma 4.1 (Regularized IHT step progress). Let $f \in \mathbb{R}^n \rightarrow \mathbb{R}$ be a convex function that is β -smooth and α -strongly convex, $\kappa = \beta/\alpha$ be its condition number, and \mathbf{x}^* be any s -sparse solution.

Given any s' -sparse solution $\mathbf{x} \in \mathbb{R}^n$ where

$$s' \geq (128\kappa + 2)s$$

and a weight vector $\mathbf{w} \in (\{0\} \cup [1/2, 1])^n$ such that $\|\mathbf{w}\|_1 \geq n - s'/2$, we make the following update:

$$\mathbf{x}' = H_{s'}((\mathbf{1} - 0.5\mathbf{w})\mathbf{x} - (2\beta)^{-1} \nabla f(\mathbf{x})).$$

Then, at least one of the following two conditions holds:

- Updating \mathbf{x} makes regularized progress:

$$g(\mathbf{x}') \leq g(\mathbf{x}) - (16\kappa)^{-1}(g(\mathbf{x}) - f(\mathbf{x}^*)),$$

where

$$g(\mathbf{x}) := f(\mathbf{x}) + (\beta/2) \|\mathbf{x}\|_{w,2}^2$$

is the ℓ_2 -regularized version of f with weights given by w . Note: The regularized progress statement is true as long as \mathbf{x} is suboptimal, i.e. $g(\mathbf{x}) > f(\mathbf{x}^*)$. Otherwise, we just have $g(\mathbf{x}') \leq g(\mathbf{x})$.

- \mathbf{x} is significantly correlated to the optimal support $S^* := \text{supp}(\mathbf{x}^*)$:

$$\|\mathbf{x}_{S^*}\|_{w^2,2}^2 \geq (4\kappa + 6)^{-1} \|\mathbf{x}\|_{w,2}^2,$$

and the regularization term restricted to S^* is non-negligible:

$$(\beta/2) \|\mathbf{x}_{S^*}\|_{w^2,2}^2 \geq (8\kappa + 8)^{-1} (g(\mathbf{x}) - f(\mathbf{x}^*)).$$

Comparison to ARHT. The ARHT algorithm of (Axiotis & Sviridenko, 2021b) is also able to achieve a sparsity bound of $O(s\kappa)$. However, their algorithm is not practically desirable for a variety of reasons.

- First of all, it follows the OMP (more accurately, OMP with Removals) paradigm, which makes local changes to the support of the solution by inserting or removing a single element of the support, and then *fully re-optimizing* the function on its restriction to this support. Even though the support will generally be very small compared to the ambient dimension n , this is still a significant runtime overhead. In contrast, regularized IHT does not require re-optimization.

Additionally, the fact that in the ARHT only one new element is added at a time leads to an iteration count that scales with $s\kappa$, instead of the κ of regularized IHT. This is a significant speedup, since both algorithms have to evaluate the gradient in each iteration. Therefore, regularized IHT will require $O(s)$ times fewer gradient evaluations.

- When faced with the non-progress condition, in which the regularized function value does not decrease sufficiently, ARHT moves by selecting a random index i with probability proportional to x_i^2 , and proceeds to *unregularize* this element, i.e. remove it from the sum of regularization terms. Instead, our algorithm is completely deterministic. This is achieved by allowing a *weighted* regularization term, and gradually reducing the regularization weights instead of dropping terms.

- ARHT requires knowledge of the optimal function value $f(\mathbf{x}^*)$. The reason is that in each iteration they need to gauge whether enough progress was made in reducing the value of the regularized function g , compared to how far it is from the optimal function value. If so, they would perform the unregularization step. In contrast, our analysis does not require these two cases (updates to \mathbf{x} or w) to be exclusive, and in fact simultaneously updates both, regardless of how much progress was made in g . Thus, our algorithm avoids the expensive overhead of an outer binary search over the optimal value $f(\mathbf{x}^*)$.

For all these reasons, as well as its striking simplicity, we believe that regularized IHT can prove to be a useful practical sparse optimization tool.

5. Low Rank Optimization Using Regularized Local Search

In this section we present a regularized local search algorithm for low rank optimization of convex functions, that returns an ε -optimal solution with rank $O\left(r\left(\kappa + \log \frac{f(\mathbf{O}) - f(\mathbf{A}^*)}{\varepsilon}\right)\right)$, where r is the target rank. The algorithm is based on the Local Search algorithm of (Axiotis & Sviridenko, 2021a), but also uses adaptive regularization, which leads to a lot new technical hurdles that are addressed in the analysis. This is presented in Theorem 1.2 and proved in Appendix D.

Theorem 1.2 (Adaptive Regularization for Low Rank Optimization). Let $f \in \mathbb{R}^{m \times n} \rightarrow \mathbb{R}$ be a convex function with condition number κ and consider the low rank minimization problem

$$\min_{\text{rank}(\mathbf{A}) \leq r} f(\mathbf{A}). \quad (8)$$

For any error parameter $\varepsilon > 0$, there exists a polynomial time algorithm that returns a matrix \mathbf{A} with $\text{rank}(\mathbf{A}) \leq O\left(r\left(\kappa + \log \frac{f(\mathbf{O}) - f(\mathbf{A}^*)}{\varepsilon}\right)\right)$ and $f(\mathbf{A}) \leq f(\mathbf{A}^*) + \varepsilon$, where \mathbf{O} is the all-zero matrix and \mathbf{A}^* is any rank- r matrix.

Discussion about ε dependence. Some of the technical issues in the rank case have to do with operator *non-commutativity* and thus pose no issue in the sparsity case. In particular, the extra $\log \frac{f(\mathbf{O}) - f(\mathbf{A}^*)}{\varepsilon}$ dependence in the rank comes exactly because of these issues. However, we think that it should be possible to completely remove this dependence in the future by a more careful analysis.

Discussion about computational efficiency. We note that the goal of this section is to show an improved rank bound, and not to argue about the computational efficiency

of such an algorithm. It might be possible to derive an efficient algorithm by transforming the proof in Theorem 1.2 into a proof for a matrix IHT algorithm, which might be significantly more efficient, as it will not require solving linear systems in each iteration. Still, there are a lot of remaining issues to be tackled, as currently the algorithm requires computing multiple singular value decompositions and orthogonal projections in each iteration. Therefore working on a computationally efficient algorithm that can guarantee a rank of $O(r\kappa)$ is a very interesting direction for future research.

Matrix regularizer Getting back into the main ingredients of Theorem 1.2, we describe the choice of our regularizer. As we are working over general rectangular matrices, we use *two* regularizers, one for the left singular vectors and one for the right singular vectors of \mathbf{A} . Concretely, given two weight matrices \mathbf{Y} , \mathbf{W} such that $\mathbf{O} \preceq \mathbf{Y} \preceq \mathbf{I}$, $\mathbf{O} \preceq \mathbf{W} \preceq \mathbf{I}$, we define

$$\Phi(\mathbf{A}) = (\beta/4) \left(\langle \mathbf{W}, \mathbf{A}\mathbf{A}^\top \rangle + \langle \mathbf{Y}, \mathbf{A}^\top \mathbf{A} \rangle \right),$$

where β is a bound on the smoothness of f . The gradient of the regularized function is

$$\nabla g(\mathbf{A}) = \nabla f(\mathbf{A}) + (\beta/2) (\mathbf{W}\mathbf{A} + \mathbf{A}\mathbf{Y}),$$

and the new solution $\bar{\mathbf{A}}$ is defined as

$$\bar{\mathbf{A}} = H_{s^{t-1}}(\mathbf{A}) - \eta H_1(\nabla g(\mathbf{A})),$$

where we remind that the thresholding operator $H_r : \mathbb{R}^{m \times n} \rightarrow \mathbb{R}^{m \times n}$ that is used in the algorithm returns the top r components of the singular value decomposition of a matrix, i.e. given $\mathbf{M} = \sum_{i=1}^k \lambda_i \mathbf{u}_i \mathbf{v}_i^\top$, where $\lambda_1 \geq \dots \geq \lambda_k$

are the singular values and $r \leq k$, $H_r(\mathbf{M}) = \sum_{i=1}^r \lambda_i \mathbf{u}_i \mathbf{v}_i^\top$.

In other words, we drop the bottom rank-1 component of \mathbf{A} and add the top rank-1 component of the gradient.

After taking a step, we re-optimize over matrices with the current left and right singular space, also known as performing a fully corrective step, as in (Shalev-Shwartz et al., 2011; Axiotis & Sviridenko, 2021a). To do this, we first compute the SVD $\mathbf{U}\mathbf{\Sigma}\mathbf{V}^\top$ of $\bar{\mathbf{A}}$ and then solve the optimization problem $\min_{\mathbf{A}=\mathbf{U}\mathbf{X}\mathbf{V}^\top} g^t(\mathbf{A})$. For simplicity we assume that this optimization problem can be solved exactly, but the analysis can be modified to account for the case when we have an approximate solution and we are only given a bound on the norm of the gradient (projected onto the relevant subspace), i.e. $\|\mathbf{\Pi}_{\text{im}(\mathbf{U})} \nabla g^t(\mathbf{A}) \mathbf{\Pi}_{\text{im}(\mathbf{V})}\|_F$.

Whenever there is not enough progress, we make the follow-

ing updates on the weight matrices \mathbf{W} and \mathbf{Y} :

$$\begin{aligned} \mathbf{W}' &= \mathbf{W} - \mathbf{W}\mathbf{A}\mathbf{A}^\top \mathbf{W} / \langle \mathbf{W}, \mathbf{A}\mathbf{A}^\top \rangle \\ \mathbf{Y}' &= \mathbf{Y} - \mathbf{Y}\mathbf{A}^\top \mathbf{A} \mathbf{Y} / \langle \mathbf{Y}, \mathbf{A}^\top \mathbf{A} \rangle. \end{aligned}$$

The full algorithm description is in Algorithm 2, Appendix D. In the algorithm description we assume that $f(\mathbf{A}^*)$ is known. This assumption can be removed by performing binary search over this value, as in (Axiotis & Sviridenko, 2021b).

6. Experiments

Introduction. In this section we present numerical experiments in order to compare the performance of IHT and regularized IHT (Algorithm 1) in training sparse linear models. In particular, we will look at the tasks of linear regression and logistic regression using both real and synthetic data. In the former, we are given a matrix $\mathbf{A} \in \mathbb{R}^{m \times n}$, where each row represents an example and each column a feature, and a vector $\mathbf{b} \in \mathbb{R}^m$ that represents the ground truth outputs, and our objective is to minimize the ℓ_2 loss

$$(1/2) \|\mathbf{A}\mathbf{x} - \mathbf{b}\|_2^2.$$

In logistic regression, \mathbf{b} has binary instead of real entries, and our objective is to minimize the logistic loss

$$-\underbrace{\sum_{i=1}^m (b_i \log \sigma(\mathbf{A}\mathbf{x})_i + (1 - b_i) \log(1 - \sigma(\mathbf{A}\mathbf{x})_i))}_{l(\mathbf{x})},$$

where $\sigma(z) = (1 + e^{-z})^{-1}$ is the sigmoid function. As is common, we look at the *regularized* logistic regression objective:

$$l(\mathbf{x}) + (\rho/2) \|\mathbf{x}\|_2^2,$$

for some $\rho > 0$. For our experiments we use $\rho = 0.1$.

Preprocessing and choice of parameters. The only preprocessing we perform is to center the columns of \mathbf{A} , i.e. we subtract the mean of each column from each entry of the column, and then scale the columns to unit ℓ_2 norm. This ensures that for any sparsity parameter $s' \in [n]$, the function f is s' -smooth when restricted to s' -sparse directions, or in other words the s' -restricted smoothness constant of f is at most s' . Thus we set our smoothness estimate to $\beta := s'$. Our smoothness estimate β influences the (regularized) IHT algorithm in two ways. First, as the step size of the algorithm is given by $1/\beta$, a value of β that is too large can slow down the algorithm, or even get it stuck to a local minimum. Second, the strength of the regularization term in regularized IHT should be close to the $(s + s')$ -restricted smoothness constant, as shown in the analysis of Theorem 1.1.

Even though having a perfectly accurate estimate of the smoothness constant is not necessary, a more accurate estimate improves the performance of the algorithm. In fact, the estimate $1/s'$ for the step size is generally too conservative. When used in practice, one should either tune this parameter or use a variable/adaptive step size to achieve the best results.

For the weight step size of regularized IHT, we set the weight step size to $c = s'/T$, but we also experiment with how changing c affects the performance of the algorithm. The downside of this setting is that it requires knowing the number of iterations a priori. However, in practice one could tune c and then run the algorithm for $O(s'/c)$ iterations. Note that ideally, based on the theoretical analysis, T would be proportional to the restricted condition number of f , however this quantity is hard to compute in general. Another idea to avoid this in practice could be to let c be a variable step size.

Implementation. Both the IHT and regularized IHT algorithms are incredibly simple, and can be described in a few lines of python code, as can be seen in Figure 2, Appendix A. Note that in comparison to Algorithm 1 we do not perform the conditional assignment. All the experiments were run on a single 2.6GHz Intel Core i7 core of a 2019 MacBook Pro with 16GB DDR4 RAM using Python 3.9.10.

6.1. Real data

We first experiment with real data, specifically the *year* regression dataset from UCI (Dua & Graff, 2017) and the *rcv1* binary classification dataset (Lewis et al., 2004), which have been previously used in the literature. Performance on other datasets was similar. In Figure 1 (a)-(b) we have a comparison between the error of the solution returned by IHT and regularized IHT for a fixed sparsity level. Specifically, if we let \mathbf{x}^{**} be the (dense) global minimizer of f , we plot the logarithm of the (normalized) *excess loss* $(f(\mathbf{x}) - f(\mathbf{x}^{**}))/f(\mathbf{0})$ against the number of iterations. Note that $f(\mathbf{x}^{**})$ will typically be considerably lower than the loss of the sparse optimum $f(\mathbf{x}^*)$. In order to make a fair comparison, for each algorithm we pick the best fixed step size of the form $2^i/s$ for integer $i \geq 0$, where s is the fixed sparsity level. The best step sizes of IHT and regularized IHT end up being $2/s, 4/s$ respectively for the linear regression example, and $8/s, 16/s$ respectively for the logistic regression example.

We notice that initially regularized IHT has a much higher error than IHT, but after some iterations it is lower than IHT. This phenomenon is to be expected, because the algorithm runs on a regularized function, and so tries to keep not just $f(\mathbf{x})$ but also $\|\mathbf{x}\|_2^2$ small. After some iterations, when the algorithm has learnt regularization weights that are closer to

the optimal ones, it converges to sparser solutions than IHT (equivalently, lower error solutions with the same sparsity, which is what is shown in the plot).

In Figure 1 (c) we compare IHT and regularized IHT for different sparsity levels on the year dataset. If e_1 and e_2 are the excess errors of IHT and regularized IHT respectively, we plot e_2/e_1 , which is the relative excess error of regularized IHT with respect to that of IHT. We notice a reduction of up to 40% on the excess error. In Figure 1 (d) we examine the effect of the choice of the weight step size c . We conclude that c can give a tradeoff between runtime and accuracy, as setting it to a large value will lead to faster weight decay and thus resemble IHT, while a small value of c will lead to slow weight decrease, which will lead to more iterations but also potentially recover an improved solution. Here we can see an interesting tradeoff between the number of iterations and the error of the solution that is eventually returned. In particular, the *larger* c is, the *faster* the degradation of regularization weights. Thus, for $c \rightarrow \infty$, the algorithm tends to be the same as IHT. On the other hand, with *smaller* values of c , one can get an improved error rate, but at the cost of a larger number of iterations. This is because the regularization weights decrease slowly, and so in the early iterations of the algorithm the regularization term will account for a significant fraction of the objective function value.

6.2. Synthetic data

We now turn to synthetically generated linear regression instances. The first result presented in Figure 1 (e) is the hard IHT instance that we derived in our lower bound in Appendix E. This experiment shows that there exist examples where, with bad initialization, IHT cannot decrease the objective at all (i.e. is stuck at a local minimum), while regularized IHT with the same initialization manages to reduce the loss by more than 70%.

The second result is a result in the well known setting of sparse signal recovery from linear measurements. We generate a matrix \mathbf{A} with entries that are sampled i.i.d. from the standard normal distribution, an s -sparse signal \mathbf{x} again with entries sampled i.i.d. from the standard normal distribution, and an observed vector $\mathbf{b} := \mathbf{A}\mathbf{x}$. The goal is to recover \mathbf{x} by minimizing the objective

$$f(\mathbf{x}) = (1/2) \|\mathbf{A}\mathbf{x} - \mathbf{b}\|_2^2.$$

In Figure 1 (f), we plot the normalized value of this objective, after running both IHT and regularized IHT for the same number of iterations. Here we pick the best step size per instance, starting from $\eta = 1/s$ and increasing in multiples of 1.2. Also, for each fixed value of s and algorithm, we run the experiments 20 times in order to account for the variance. The results show a superiority in the performance of regularized IHT for the sparse signal recovery task.

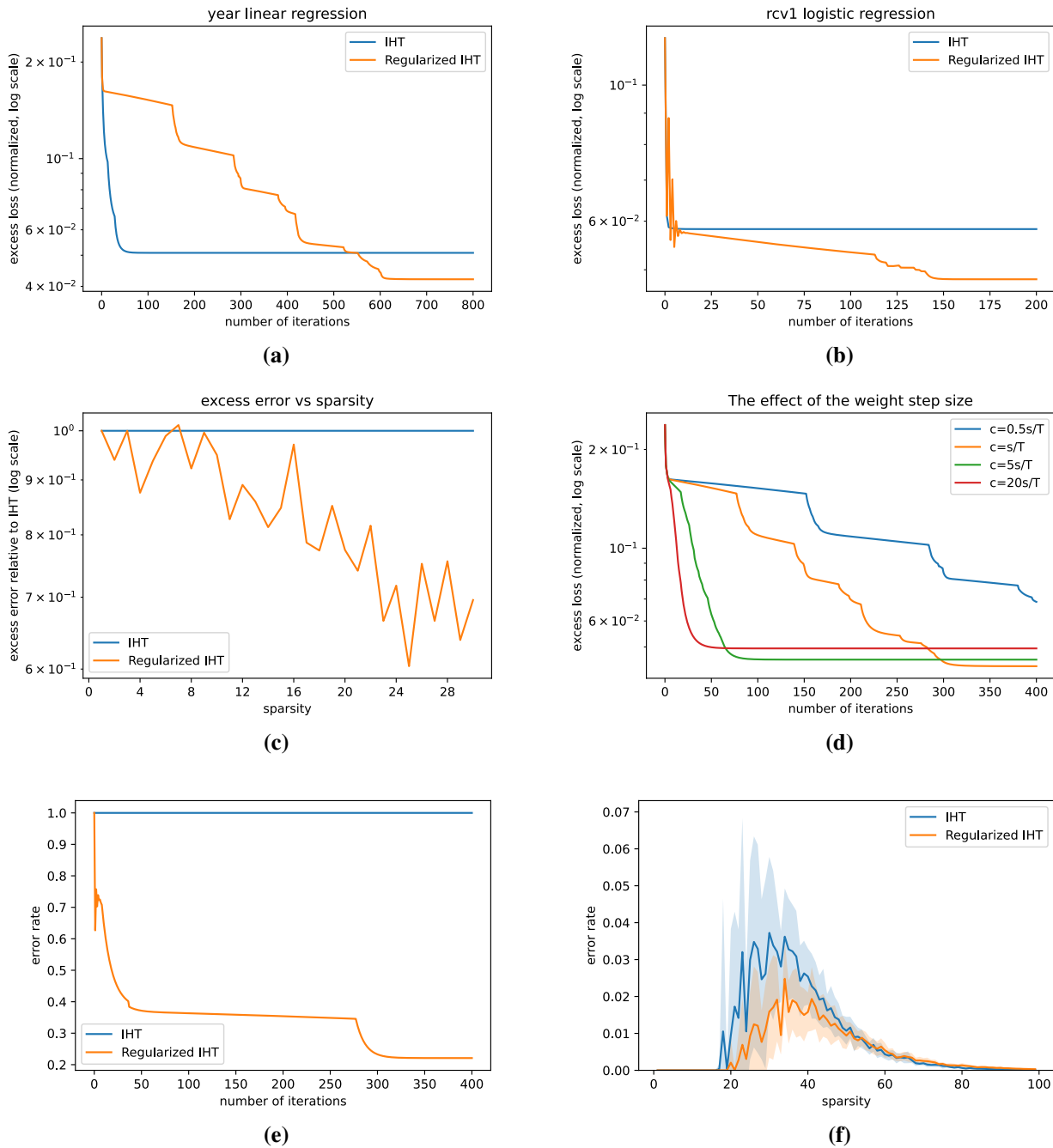


Figure 1. (a)-(b): IHT vs Regularized IHT performance on the year and rcv1 datasets with fixed sparsity levels $s = 11$ and $s = 10$ respectively. On the x axis we have number of iterations and on the y axis we have the normalized excess loss (compared to the dense global optimum), in log scale. The excess loss of regularized IHT is less than that of IHT, specifically 17.3% and 17.2% respectively less in the two experiments. (c): Excess error of regularized IHT relative to IHT in the year dataset, where sparsity values range from 1 to 30. Both algorithms are run for $T = 800$ iterations. (d): Error rate vs number of iterations of regularized IHT on the year dataset with fixed sparsity $s = 11$ and step size $\eta = 4/s$, using different values for the weight step size c . (e): A demonstration of a 80% decrease in loss by using regularized IHT instead of IHT on the hard instance for IHT presented in Section E. We have generated the data with a condition number of $\kappa = 20$, and a planted sparse solution with sparsity $s = 2$. The dimension is $n = 842$. It can be observed that, for the given initialization vector, IHT never makes any progress on decreasing the error. In contrast, regularized IHT is able to decrease it by almost a factor of 5. (f): Sparse signal recovery, where \mathbf{A} is an 100×800 measurement matrix, the sparsity level ranges from 1 to 100, and each algorithm is run for 240 iterations. Bands of 1 standard error are shown, after running each data point 20 times independently.

References

- Allen-Zhu, Z., Hazan, E., Hu, W., and Li, Y. Linear convergence of a frank-wolfe type algorithm over trace-norm balls. *arXiv preprint arXiv:1708.02105*, 2017.
- Axiotis, K. and Sviridenko, M. Local search algorithms for rank-constrained convex optimization. In *International Conference on Learning Representations*, 2021a.
- Axiotis, K. and Sviridenko, M. Sparse convex optimization via adaptively regularized hard thresholding. *Journal of Machine Learning Research*, 22:1–47, 2021b.
- Blumensath, T. and Davies, M. E. Iterative hard thresholding for compressed sensing. *Applied and computational harmonic analysis*, 27(3):265–274, 2009.
- Boche, H., Calderbank, R., Kutyniok, G., and Vybíral, J. A survey of compressed sensing. In *Compressed sensing and its applications*, pp. 1–39. Springer, 2015.
- Candes, E. J. The restricted isometry property and its implications for compressed sensing. *Comptes rendus mathématique*, 346(9-10):589–592, 2008.
- Donoho, D. L. and Elad, M. Optimally sparse representation in general (nonorthogonal) dictionaries via ℓ_1 minimization. *Proceedings of the National Academy of Sciences*, 100(5):2197–2202, 2003.
- Donoho, D. L., Maleki, A., and Montanari, A. Message-passing algorithms for compressed sensing. *Proceedings of the National Academy of Sciences*, 106(45):18914–18919, 2009.
- Dua, D. and Graff, C. UCI machine learning repository, 2017. URL <http://archive.ics.uci.edu/ml>.
- Duchi, J., Hazan, E., and Singer, Y. Adaptive subgradient methods for online learning and stochastic optimization. *Journal of machine learning research*, 12(7), 2011.
- Foster, D., Karloff, H., and Thaler, J. Variable selection is hard. In *Conference on Learning Theory*, pp. 696–709, 2015.
- Foucart, S. A note on guaranteed sparse recovery via ℓ_1 -minimization. *Applied and Computational Harmonic Analysis*, 29(1):97–103, 2010.
- Foucart, S. Hard thresholding pursuit: an algorithm for compressive sensing. *SIAM Journal on Numerical Analysis*, 49(6):2543–2563, 2011.
- Foucart, S. and Rauhut, H. A mathematical introduction to compressive sensing. *Bull. Am. Math.*, 54:151–165, 2017.
- Frank, M., Wolfe, P., et al. An algorithm for quadratic programming. *Naval research logistics quarterly*, 3(1-2): 95–110, 1956.
- Horn, R. A. and Johnson, C. Topics in matrix analysis, 1991.
- Jaggi, M. Revisiting frank-wolfe: Projection-free sparse convex optimization. In *International Conference on Machine Learning*, pp. 427–435. PMLR, 2013.
- Jain, P., Tewari, A., and Dhillon, I. S. Orthogonal matching pursuit with replacement. In *Advances in neural information processing systems*, pp. 1215–1223, 2011.
- Jain, P., Tewari, A., and Kar, P. On iterative hard thresholding methods for high-dimensional m -estimation. In *Advances in Neural Information Processing Systems*, pp. 685–693, 2014.
- Jain, P., Tewari, A., and Dhillon, I. S. Partial hard thresholding. *IEEE Transactions on Information Theory*, 63(5): 3029–3038, 2017.
- Lewis, D. D., Yang, Y., Russell-Rose, T., and Li, F. Rcv1: A new benchmark collection for text categorization research. *Journal of machine learning research*, 5(Apr):361–397, 2004.
- Liberty, E. and Sviridenko, M. Greedy minimization of weakly supermodular set functions. In *Approximation, Randomization, and Combinatorial Optimization. Algorithms and Techniques (APPROX/RANDOM 2017)*. Schloss Dagstuhl-Leibniz-Zentrum fuer Informatik, 2017.
- Liu, H. and Foygel Barber, R. Between hard and soft thresholding: optimal iterative thresholding algorithms. *Information and Inference: A Journal of the IMA*, 9(4): 899–933, 2020.
- Natarajan, B. K. Sparse approximate solutions to linear systems. *SIAM journal on computing*, 24(2):227–234, 1995.
- Needell, D. and Tropp, J. A. Cosamp: Iterative signal recovery from incomplete and inaccurate samples. *Applied and computational harmonic analysis*, 26(3):301–321, 2009.
- Pati, Y. C., Rezaiifar, R., and Krishnaprasad, P. S. Orthogonal matching pursuit: Recursive function approximation with applications to wavelet decomposition. In *Proceedings of 27th Asilomar conference on signals, systems and computers*, pp. 40–44. IEEE, 1993.
- Peste, A., Iofinova, E., Vladu, A., and Alistarh, D. Ac/dc: Alternating compressed/decompressed training of deep neural networks. In *Thirty-Fifth Conference on Neural Information Processing Systems*, 2021.

Schmidt, L. *Algorithms above the noise floor*. PhD thesis, Massachusetts Institute of Technology, Cambridge, USA, 2018. URL <http://hdl.handle.net/1721.1/118098>.

Shalev-Shwartz, S., Srebro, N., and Zhang, T. Trading accuracy for sparsity in optimization problems with sparsity constraints. *SIAM Journal on Optimization*, 20(6): 2807–2832, 2010.

Shalev-Shwartz, S., Gonen, A., and Shamir, O. Large-scale convex minimization with a low-rank constraint. In *Proceedings of the 28th International Conference on International Conference on Machine Learning*, pp. 329–336, 2011.

Shen, J. and Li, P. On the iteration complexity of support recovery via hard thresholding pursuit. In *Proceedings of the 34th International Conference on Machine Learning-Volume 70*, pp. 3115–3124. JMLR. org, 2017a.

Shen, J. and Li, P. Partial hard thresholding: Towards a principled analysis of support recovery. In *Advances in Neural Information Processing Systems*, pp. 3124–3134, 2017b.

Tibshirani, R. Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society: Series B (Methodological)*, 58(1):267–288, 1996.

Yuan, X., Li, P., and Zhang, T. Exact recovery of hard thresholding pursuit. In *Advances in Neural Information Processing Systems*, pp. 3558–3566, 2016.

A. Python implementation

In Figure 2 we have python implementations of the IHT and regularized IHT algorithms that we use for our experiments. As can be seen, both implementations are pretty short.

B. Proof of Theorem 1.1

Theorem 1.1 (Regularized IHT). Let $f \in \mathbb{R}^n \rightarrow \mathbb{R}$ be a convex function that is β -smooth and α -strongly convex, with condition number $\kappa = \beta/\alpha$, and \mathbf{x}^* be an (unknown) s -sparse solution with support S^* . Then, running Algorithm 1 with $\eta = (2\beta)^{-1}$ and $c = s'/(4T)$ for

$$T = O\left(\kappa \log \frac{f(\mathbf{x}^0) + (\beta/2) \|\mathbf{x}^0\|_2^2 - f(\mathbf{x}^*)}{\varepsilon}\right)$$

iterations starting from an arbitrary $s' = O(s\kappa)$ -sparse solution \mathbf{x}^0 , the algorithm returns an s' -sparse solution \mathbf{x}^T such that $f(\mathbf{x}^T) \leq f(\mathbf{x}^*) + \varepsilon$. Furthermore, each iteration requires $O(1)$ evaluations of f , ∇f , and $O(n)$ additional time.

Proof. We repeatedly apply Lemma 4.1 for

$$T = 64(\kappa + 1) \log \frac{f(\mathbf{x}^0) + (\beta/2) \|\mathbf{x}^0\|_2^2 - f(\mathbf{x}^*)}{\varepsilon}$$

iterations. We define the regularized function

$$g^t(\mathbf{x}) := f(\mathbf{x}) + (\beta/2) \|\mathbf{x}\|_{\mathbf{w}^t, 2}^2,$$

where \mathbf{w}^t are the weights before iteration $t \in [0, T - 1]$. Specifically, for each t we apply Lemma 4.1 on the current solution \mathbf{x}^t and obtain the solution \mathbf{x}^{t+1} .

Before moving forward, we give an intuitive summary of the proof and the role of Lemma 4.1. As long as IHT makes “sufficient” progress on the regularized function g^t , this is satisfactory for the original function f as well, because $f(\mathbf{x}) \leq g^t(\mathbf{x})$ for all \mathbf{x} . This is the case of the first bullet of Lemma 4.1. If it stops making sufficient progress, this means we are at an (approximate) sparse optimum for g^t , although it is not necessarily a good sparse solution for f , which is the objective we are aiming to minimize. This is where the second bullet of Lemma 4.1 comes in, which gives necessary conditions for the above non-progress phenomenon (in other words, a partial characterization of the local minima encountered when running IHT on a regularized function). Specifically, the following condition is central to our approach:

$$\|\mathbf{x}_{S^*}^t\|_{(\mathbf{w}^t)^2, 2}^2 \geq (4\kappa + 6)^{-1} \|\mathbf{x}^t\|_{\mathbf{w}^t, 2}^2.$$

We use this condition in the second part of the proof (after Case 2) to motivate a weight update from \mathbf{w}^t to \mathbf{w}^{t+1} , and

```

import numpy as np

def IHT(n, s):
    x = np.zeros(n)
    for _ in range(T):
        x_new = x - eta * grad(x)
        x_new[np.argsort(np.abs(x_new))[:s]] = 0
        x = x_new
    return x

def RegIHT(n, s):
    x, w = np.zeros(n), np.ones(n)
    for _ in range(T):
        x_new = (1 - 0.5 * w) * x - 0.5 * eta * grad(x)
        x_new[np.argsort(np.abs(x_new))[:s]] = 0
        reg = np.sum(w * x**2)
        if reg != 0:
            w = w * (1 - c * w * x**2 / reg)
            w[w <= round_th] = 0
        x = x_new
    return x
    
```

Figure 2. Our python implementations of IHT, RegIHT, where grad is the gradient function, n is the total number of features, s is the desired sparsity level, η is the step size, c is the weight step size, and round_th is the weight rounding threshold, which we set to 0.5. Note that $\text{grad}(x) = \text{np.dot}(A.T, \text{np.dot}(A, x) - b)$ for linear regression and $\text{grad}(x) = \text{np.dot}(A.T, \text{expit}(\text{np.dot}(A, x)) - b)$ for logistic regression, where expit is the sigmoid function.

show that, exactly because of this condition, a lot of the weight decrease is concentrated inside the optimal support S^* . As the total weight decrease in S^* is bounded by s , this gives a bound on the total number of iterations with insufficient decrease of g^t . If not for this condition, we would not be able to bound the number of such iterations and would have potentially remained forever stuck at a local minimum.

Now we are ready to move to the technical proof. In order to make sure that $g^{t+1}(\mathbf{x}^{t+1}) \leq g^{t+1}(\mathbf{x}^t)$, we revert to the previous solution if the one returned by Lemma 4.1 has a larger value of g^{t+1} . This is exactly what the conditional in Algorithm 1 is for. The property that $g^{t+1}(\mathbf{x}^{t+1}) \leq g^{t+1}(\mathbf{x}^t)$ is only used in the very last part of the proof.

Let us assume that $g^t(\mathbf{x}^t) > f(\mathbf{x}^*)$ at all times, as otherwise the statement holds by the fact that $g^t(\mathbf{x}^t)$ is non-increasing as a function of t and upper bounds $f(\mathbf{x}^t)$ for all t . We have $g^{t+1}(\mathbf{x}^{t+1}) \leq g^t(\mathbf{x}^t)$ by the fact that $\mathbf{w}^{t+1} \leq \mathbf{w}^t$ and $g^t(\mathbf{x}^{t+1}) \leq g^t(\mathbf{x}^t)$ by the guarantees of Lemma 4.1.

If the first bullet of Lemma 4.1 holds, we have that the value of g decreases considerably on iteration t , i.e.

$$\begin{aligned}
 & g^{t+1}(\mathbf{x}^{t+1}) \\
 & \leq g^t(\mathbf{x}^{t+1}) \\
 & \leq g^t(\mathbf{x}^t) - (16\kappa)^{-1}(g^t(\mathbf{x}^t) - f(\mathbf{x}^*)).
 \end{aligned}$$

Let us call these iterations *progress iterations*, and the other ones (where the second bullet of Lemma 4.1 holds) *weight*

iterations. Now, since $g^t(\mathbf{x}^t)$ is non-increasing as a function of t , after $16\kappa \log \frac{g^0(\mathbf{x}^0) - f(\mathbf{x}^*)}{\varepsilon}$ progress iterations we will have

$$f(\mathbf{x}^t) \leq g^t(\mathbf{x}^t) \leq f(\mathbf{x}^*) + \varepsilon,$$

and so we will be done. From now on let us assume this is not the case, so there are at least

$$T - 16\kappa \log \frac{g^0(\mathbf{x}^0) - f(\mathbf{x}^*)}{\varepsilon} \geq 3T/4$$

weight iterations.

We remind that in each weight iteration, we have

$$\begin{aligned}
 & \|\mathbf{x}_{S^*}^t\|_{(\mathbf{w}^t)^2, 2}^2 \geq (4\kappa + 6)^{-1} \|\mathbf{x}^t\|_{\mathbf{w}^t, 2}^2 \quad (9) \\
 & (\beta/2) \|\mathbf{x}_{S^*}^t\|_{(\mathbf{w}^t)^2, 2}^2 \geq (8\kappa + 8)^{-1} (g^t(\mathbf{x}^t) - f(\mathbf{x}^*)). \quad (10)
 \end{aligned}$$

In words, (9) roughly implies that at least an $\Omega(1/\kappa)$ fraction of the mass of $(\mathbf{w}^t \mathbf{x}^t)^2$ lies inside S^* . Therefore, if we decrease \mathbf{w}^t by a quantity proportional to $(\mathbf{w}^t \mathbf{x}^t)^2$, the total sum of weights will decrease at most $O(\kappa)$ times faster than the sum of weights inside S^* . As the latter quantity can only decrease by s overall, the total decrease of weights will be $O(s\kappa)$.

Concretely, after each iteration we update the regularization weights as follows:

$$\mathbf{w}^{t+1} = \left(\mathbf{w}^t - c \cdot (\mathbf{w}^t \mathbf{x}^t)^2 / \|\mathbf{x}^t\|_{\mathbf{w}^t, 2}^2 \right)_{\geq 1/2},$$

for some $c > 0$ to be determined later. First of all, note that the weights are non-increasing. Now, if not for the thresholding operation, it is easy to see that the total weight decrease is at most c . The thresholding operation can only double this weight decrease to $2c$. Concretely, for all t we define a vector \bar{w}^t such that

$$\bar{w}_i^t = \begin{cases} w_i^t & \text{if } w_i^t \geq 1/2 \\ 1/2 & \text{if } w_i^t = 0. \end{cases}$$

Clearly, $\frac{1}{2}(1 - w^t) \leq 1 - \bar{w}^t \leq 1 - w^t$. Now, we have

$$\begin{aligned} & \|\bar{w}^t\|_1 - \|\bar{w}^{t+1}\|_1 \\ & \leq c \|\mathbf{x}^t\|_{(w^t)^2,2}^2 / \|\mathbf{x}^t\|_{w^t,2}^2 \\ & \leq c, \end{aligned}$$

and, summing up for all t we get

$$\|\mathbf{1} - \bar{w}^T\|_1 = \|\bar{w}^0\|_1 - \|\bar{w}^T\|_1 \leq cT.$$

Therefore,

$$\|\mathbf{1} - w^T\|_1 \leq 2\|\mathbf{1} - \bar{w}^T\|_1 \leq 2cT,$$

and so $\|\mathbf{w}^T\|_1 \geq n - 2cT$.

Therefore, the condition $\|\mathbf{w}^t\|_1 \geq n - s'/2$ of Lemma 4.1 is satisfied for all t as long as $c \leq s'/(4T)$. In order to bound the number of iterations, we distinguish two cases for the sum of weights inside S^* .

Case 1: The sum of weights inside S^* decreases by $\geq 4s/T$.

This case cannot happen more than $T/4$ times since the sum of weights inside S^* can only decrease by s in total. Therefore, case 2 below happens at least $T/2$ times.

Case 2: The sum of weights inside S^* decreases by $< 4s/T$.

Note that the decrease in the sum of weights in S^* is exactly equal to

$$\sum_{i \in S^*} \begin{cases} c \cdot (w_i^t x_i^t)^2 / \|\mathbf{x}^t\|_{w^t,2}^2 & \text{if this is } \leq w_i^t - 1/2 \\ w_i^t & \text{otherwise.} \end{cases}$$

Let T^* be the set of indices $i \in S^*$ for which the second case is true, i.e.

$$c \cdot (w_i^t x_i^t)^2 / \|\mathbf{x}^t\|_{w^t,2}^2 > w_i^t - 1/2.$$

The total weight decrease from elements in $S^* \setminus T^*$ is then

$$\begin{aligned} & \sum_{i \in S^* \setminus T^*} c \cdot (w_i^t x_i^t)^2 / \|\mathbf{x}^t\|_{w^t,2}^2 \\ & = c \|\mathbf{x}_{S^* \setminus T^*}^t\|_{(w^t)^2,2}^2 / \|\mathbf{x}^t\|_{w^t,2}^2 \\ & \geq \frac{c}{4\kappa + 6} \|\mathbf{x}_{S^* \setminus T^*}^t\|_{(w^t)^2,2}^2 / \|\mathbf{x}_{S^*}^t\|_{(w^t)^2,2}^2, \end{aligned}$$

where we used (9). As we have assumed that this decrease is less than $4s/T$, we have that

$$\begin{aligned} \|\mathbf{x}_{T^*}^t\|_{(w^t)^2,2}^2 & = \|\mathbf{x}_{S^*}^t\|_{(w^t)^2,2}^2 - \|\mathbf{x}_{S^* \setminus T^*}^t\|_{(w^t)^2,2}^2 \\ & \geq \left(1 - \frac{4s(4\kappa + 6)}{cT}\right) \|\mathbf{x}_{S^*}^t\|_{(w^t)^2,2}^2 \quad (11) \\ & \geq (1/2) \|\mathbf{x}_{S^*}^t\|_{(w^t)^2,2}^2, \end{aligned}$$

as long as $c \geq 8s(4\kappa + 6)/T$. We can pick such a c as long as

$$8s(4\kappa + 6)/T \leq c \leq s'/(4T) \Leftrightarrow s' \geq 32(4\kappa + 6)s.$$

Now, to deal with the fact that the sum weights in T^* might not decrease sufficiently, note that all the weights in T^* are being set to 0, i.e. $w_i^{t+1} = 0$ for all $i \in T^*$. Together with (11) and (10) this means that we can make significant progress in function value. To see this, note that

$$\begin{aligned} & g^{t+1}(\mathbf{x}^{t+1}) \\ & \leq g^{t+1}(\mathbf{x}^t) \\ & \leq g^t(\mathbf{x}^t) - (\beta/2) \|\mathbf{x}_{T^*}^t\|_{w^t,2}^2 \\ & \leq g^t(\mathbf{x}^t) - (\beta/2) \|\mathbf{x}_{T^*}^t\|_{(w^t)^2,2}^2 \\ & \leq g^t(\mathbf{x}^t) - (\beta/4) \|\mathbf{x}_{S^*}^t\|_{(w^t)^2,2}^2 \\ & \leq g^t(\mathbf{x}^t) - (16\kappa + 16)^{-1} (g^t(\mathbf{x}^t) - f(\mathbf{x}^*)), \end{aligned}$$

which can happen at most

$$16(\kappa + 1) \log \frac{g^0(\mathbf{x}^0) - f(\mathbf{x}^*)}{\varepsilon} \leq T/4$$

times. □

C. Proof of Lemma 4.1

Lemma 4.1 (Regularized IHT step progress). *Let $f \in \mathbb{R}^n \rightarrow \mathbb{R}$ be a convex function that is β -smooth and α -strongly convex, $\kappa = \beta/\alpha$ be its condition number, and \mathbf{x}^* be any s -sparse solution.*

Given any s' -sparse solution $\mathbf{x} \in \mathbb{R}^n$ where

$$s' \geq (128\kappa + 2)s$$

and a weight vector $\mathbf{w} \in (\{0\} \cup [1/2, 1])^n$ such that $\|\mathbf{w}\|_1 \geq n - s'/2$, we make the following update:

$$\mathbf{x}' = H_{s'}((\mathbf{1} - 0.5\mathbf{w})\mathbf{x} - (2\beta)^{-1}\nabla f(\mathbf{x})).$$

Then, at least one of the following two conditions holds:

- Updating \mathbf{x} makes regularized progress:

$$g(\mathbf{x}') \leq g(\mathbf{x}) - (16\kappa)^{-1}(g(\mathbf{x}) - f(\mathbf{x}^*)),$$

where

$$g(\mathbf{x}) := f(\mathbf{x}) + (\beta/2) \|\mathbf{x}\|_{\mathbf{w},2}^2$$

is the ℓ_2 -regularized version of f with weights given by \mathbf{w} . Note: The regularized progress statement is true as long as \mathbf{x} is suboptimal, i.e. $g(\mathbf{x}) > f(\mathbf{x}^*)$. Otherwise, we just have $g(\mathbf{x}') \leq g(\mathbf{x})$.

- \mathbf{x} is significantly correlated to the optimal support $S^* := \text{supp}(\mathbf{x}^*)$:

$$\|\mathbf{x}_{S^*}\|_{\mathbf{w}^2,2}^2 \geq (4\kappa + 6)^{-1} \|\mathbf{x}\|_{\mathbf{w},2}^2,$$

and the regularization term restricted to S^* is non-negligible:

$$(\beta/2) \|\mathbf{x}_{S^*}\|_{\mathbf{w}^2,2}^2 \geq (8\kappa + 8)^{-1} (g(\mathbf{x}) - f(\mathbf{x}^*)).$$

Proof. By using the fact that f is β -smooth, and so g is 2β -smooth due to $\mathbf{w} \leq \mathbf{1}$, for any \mathbf{x}' we obtain

$$g(\mathbf{x}') - g(\mathbf{x}) \leq \langle \nabla g(\mathbf{x}), \mathbf{x}' - \mathbf{x} \rangle + \beta \|\mathbf{x}' - \mathbf{x}\|_2^2. \quad (12)$$

We let S be the support of \mathbf{x} and S' the support of \mathbf{x}' , i.e. the set of s' indices of the largest magnitude entries of the vector. Since $\nabla g(\mathbf{x}) = \nabla f(\mathbf{x}) + \beta \mathbf{w} \mathbf{x}$, we have

$$\bar{\mathbf{x}} = \mathbf{x} - \eta \nabla g(\mathbf{x}) = (\mathbf{1} - 0.5\mathbf{w})\mathbf{x} - \eta \nabla f(\mathbf{x}),$$

where $\eta = (2\beta)^{-1}$. We let $A = S' \setminus S$ be the newly inserted entries and $B = S \setminus S'$ be the entries that were just removed from the support. Note that

$$\begin{aligned} \mathbf{x}' &= [\mathbf{x} - \eta \nabla g(\mathbf{x})]_{S'} \\ &= \mathbf{x} - \eta \nabla_{S'} g(\mathbf{x}) - \mathbf{x}_B \\ &= \mathbf{x} - \eta \nabla_{S' \cup B} g(\mathbf{x}) - \bar{\mathbf{x}}_B. \end{aligned}$$

Using (12), we have

$$\begin{aligned} g(\mathbf{x}') - g(\mathbf{x}) &\leq \langle \nabla g(\mathbf{x}), -\eta \nabla_{S' \cup B} g(\mathbf{x}) - \bar{\mathbf{x}}_B \rangle \\ &\quad + \beta \|\eta \nabla_{S' \cup B} g(\mathbf{x}) - \bar{\mathbf{x}}_B\|_2^2 \\ &= -(4\beta)^{-1} \|\nabla_{S' \cup B} g(\mathbf{x})\|_2^2 + \beta \|\bar{\mathbf{x}}_B\|_2^2 \\ &= -\beta \|\eta \nabla_{S \cup A} g(\mathbf{x})\|_2^2 + \beta \|\bar{\mathbf{x}}_B\|_2^2 \\ &\leq -\beta \|\eta \nabla_{S \cup A'} g(\mathbf{x})\|_2^2 + \beta \|\bar{\mathbf{x}}_{B'}\|_2^2, \end{aligned} \quad (13)$$

for any two sets $A' \in [n] \setminus S$ and $B' \subseteq S$ with $|A'| = |B'|$. The latter inequality follows because of the following lemma about IHT:

Lemma C.1. Suppose that we run one step of IHT on vector \mathbf{x} supported on S for some function g , and let the updated solution vector be $\mathbf{x}' = \bar{\mathbf{x}}_{(S \cup A) \setminus B}$, where $\bar{\mathbf{x}} = \mathbf{x} - \eta \nabla g(\mathbf{x})$. Then, for any $A' \subseteq [n] \setminus S$ and $B' \subseteq S$ with $|A'| = |B'|$, we have

$$-\|\eta \nabla_{A'} g(\mathbf{x})\|_2^2 + \|\bar{\mathbf{x}}_B\|_2^2 \leq -\|\eta \nabla_{A'} g(\mathbf{x})\|_2^2 + \|\bar{\mathbf{x}}_{B'}\|_2^2. \quad (14)$$

Proof. If we denote $|A| = |B| = t$ and $|A'| = |B'| = t'$, then note that by definition of IHT, A are the t largest entries in

$$|\bar{\mathbf{x}}_{[n] \setminus S}| = \eta |\nabla_{[n] \setminus S} g(\mathbf{x})|,$$

and B are the t smallest entries in $|\bar{\mathbf{x}}_S|$. Similarly, we can assume that A' are the t' largest entries in $\eta |\nabla_{[n] \setminus S} g(\mathbf{x})|$ and B' are the t' smallest entries in $|\bar{\mathbf{x}}_S|$, since this way the right hand side of (14) takes its minimum value. If $t' = t$, we are done. We consider two cases:

1. $t' > t$: In this case we have $A' \supseteq A, B' \supseteq B$, so

$$\begin{aligned} &-\|\eta \nabla_{A'} g(\mathbf{x})\|_2^2 + \|\bar{\mathbf{x}}_{B'}\|_2^2 + \|\eta \nabla_A g(\mathbf{x})\|_2^2 - \|\bar{\mathbf{x}}_B\|_2^2 \\ &= -\|\eta \nabla_{A' \setminus A} g(\mathbf{x})\|_2^2 + \|\bar{\mathbf{x}}_{B' \setminus B}\|_2^2 \\ &= -\|\bar{\mathbf{x}}_{A' \setminus A}\|_2^2 + \|\bar{\mathbf{x}}_{B' \setminus B}\|_2^2 \\ &\geq (t' - t) \left(-\max_{i \in A' \setminus A} (\bar{x}_i)^2 + \min_{j \in B' \setminus B} (\bar{x}_j)^2 \right) \\ &\geq 0, \end{aligned}$$

where the last inequality follows since, by definition of the IHT step, $|\bar{x}_i| \leq |\bar{x}_j|$ for any $i \in A' \setminus A$ and $j \in B' \setminus B$. Otherwise, i would have taken j 's place in S' .

2. $t' < t$: In this case we have $A' \subseteq A, B' \subseteq B$. Similarly to the previous case,

$$\begin{aligned} &-\|\eta \nabla_{A'} g(\mathbf{x})\|_2^2 + \|\bar{\mathbf{x}}_{B'}\|_2^2 + \|\eta \nabla_A g(\mathbf{x})\|_2^2 - \|\bar{\mathbf{x}}_B\|_2^2 \\ &= \|\eta \nabla_{A \setminus A'} g(\mathbf{x})\|_2^2 - \|\bar{\mathbf{x}}_{B \setminus B'}\|_2^2 \\ &= \|\bar{\mathbf{x}}_{A \setminus A'}\|_2^2 - \|\bar{\mathbf{x}}_{B \setminus B'}\|_2^2 \\ &\geq (t - t') \left(\min_{i \in A \setminus A'} (\bar{x}_i)^2 - \max_{j \in B \setminus B'} (\bar{x}_j)^2 \right) \\ &\geq 0, \end{aligned}$$

where the last inequality follows since, by definition of the IHT step, $|\bar{x}_i| \geq |\bar{x}_j|$ for any $i \in A \setminus A'$ and $j \in B \setminus B'$. Otherwise i wouldn't have taken j 's place in S' .

□

Now, let us assume that the first bullet in the lemma statement is false, i.e.

$$g(\mathbf{x}') - g(\mathbf{x}) > -(16\kappa)^{-1}(g(\mathbf{x}) - f(\mathbf{x}^*)).$$

Setting $A' = B' = \emptyset$ in (13), we get that

$$g(\mathbf{x}') - g(\mathbf{x}) \leq -\beta \|\eta \nabla_S g(\mathbf{x})\|_2^2,$$

so we conclude that

$$\|\nabla_S g(\mathbf{x})\|_2^2 < \frac{1}{16\kappa\beta\eta^2} (g(\mathbf{x}) - f(\mathbf{x}^*)) = \frac{\alpha}{4} (g(\mathbf{x}) - f(\mathbf{x}^*)). \quad (15)$$

Now, we again use (13) but we set A' to be the s entries from $[n] \setminus S$ on which $\nabla g(\mathbf{x})$ has the largest magnitude, and B' to be the s entries from S on which $\bar{\mathbf{x}}$ has the smallest magnitude. Also, let R be an arbitrary subset of $S \setminus S^*$ with size $r \geq 2s$. We then have

$$\begin{aligned} g(\mathbf{x}') - g(\mathbf{x}) &\leq -(4\beta)^{-1} \|\nabla_{S \cup A'} g(\mathbf{x})\|_2^2 + \beta \|\bar{\mathbf{x}}_{B'}\|_2^2 \\ &\leq -(4\beta)^{-1} \|\nabla_{S \cup S^*} g(\mathbf{x})\|_2^2 + \beta \|\bar{\mathbf{x}}_{B'}\|_2^2 \\ &\leq -(4\beta)^{-1} \|\nabla_{S \cup S^*} g(\mathbf{x})\|_2^2 + \frac{\beta s}{|R \setminus S^*|} \|\bar{\mathbf{x}}_{R \setminus S^*}\|_2^2 \\ &\leq -(4\beta)^{-1} \|\nabla_{S \cup S^*} g(\mathbf{x})\|_2^2 + \frac{\beta s}{r-s} \|\bar{\mathbf{x}}_{R \setminus S^*}\|_2^2, \end{aligned} \quad (16)$$

where we used the fact that

$$\|\nabla_{A'} g(\mathbf{x})\|_2^2 \geq \|\nabla_{S^* \setminus S} g(\mathbf{x})\|_2^2$$

by definition of A' (and since $|S^* \setminus S| \leq s$), and the fact that, by definition of B' (and since $|R \setminus S^*| \geq 2s - s = |B'|$),

$$\frac{1}{|B'|} \|\bar{\mathbf{x}}_{B'}\|_2^2 \leq \frac{1}{|R \setminus S^*|} \|\bar{\mathbf{x}}_{R \setminus S^*}\|_2^2.$$

In fact, we will let $R = \{i \in S \mid w_i > 0\}$ be the set of elements that are being regularized. To lower bound the size r of this set, note that by the guarantee of the lemma statement,

$$\begin{aligned} n - s'/2 &\leq \|\mathbf{w}\|_1 \\ &\leq n - |\{i \in S \mid w_i = 0\}| \\ &= n - (s' - r), \end{aligned}$$

so $r \geq s'/2$. We conclude that $r \geq 2s$ since $s' \geq 4s$.

Now, because of the fact that f is α -strongly convex, we have

$$\begin{aligned} f(\mathbf{x}^*) - f(\mathbf{x}) &\geq \langle \nabla f(\mathbf{x}), \mathbf{x}^* - \mathbf{x} \rangle + (\alpha/2) \|\mathbf{x}^* - \mathbf{x}\|_2^2 \\ &= \langle \nabla g(\mathbf{x}), \mathbf{x}^* - \mathbf{x} \rangle - \beta \langle \mathbf{w}\mathbf{x}, \mathbf{x}^* - \mathbf{x} \rangle + (\alpha/2) \|\mathbf{x}^* - \mathbf{x}\|_2^2 \\ &\geq -\alpha^{-1} \|\nabla_{S \cup S^*} g(\mathbf{x})\|_2^2 \\ &\quad - \beta \langle \mathbf{w}\mathbf{x}, \mathbf{x}^* - \mathbf{x} \rangle + (\alpha/4) \|\mathbf{x}^* - \mathbf{x}\|_2^2, \end{aligned} \quad (17)$$

where we used the inequality

$$\langle \mathbf{a}, \mathbf{b} \rangle + (\alpha/4) \|\mathbf{b}\|_2^2 \geq -\alpha^{-1} \|\mathbf{a}\|_2^2.$$

By re-arranging and plugging (17) into (16), we get

$$\begin{aligned} g(\mathbf{x}') - g(\mathbf{x}) &\leq -(4\kappa)^{-1} (f(\mathbf{x}) - f(\mathbf{x}^*)) \\ &\quad - \beta \langle \mathbf{w}\mathbf{x}, \mathbf{x}^* - \mathbf{x} \rangle + (\alpha/4) \|\mathbf{x}^* - \mathbf{x}\|_2^2 \\ &\quad + \frac{\beta s}{r-s} \|\bar{\mathbf{x}}_{R \setminus S^*}\|_2^2 \\ &= -(4\kappa)^{-1} (g(\mathbf{x}) - f(\mathbf{x}^*)) \\ &\quad - \beta \langle \mathbf{w}\mathbf{x}, \mathbf{x}^* - \mathbf{x} \rangle - (\beta/2) \|\mathbf{x}\|_{\mathbf{w},2}^2 \\ &\quad + (\alpha/4) \|\mathbf{x}^* - \mathbf{x}\|_2^2 - \frac{4\kappa\beta s}{r-s} \|\bar{\mathbf{x}}_{R \setminus S^*}\|_2^2. \end{aligned} \quad (18)$$

Now, note that by definition of $\bar{\mathbf{x}}$ we have

$$\begin{aligned} \|\bar{\mathbf{x}}_{R \setminus S^*}\|_2^2 &\leq 2 \|\mathbf{x}_{R \setminus S^*}\|_2^2 + 2(2\beta)^{-2} \|\nabla_{R \setminus S^*} g(\mathbf{x})\|_2^2 \end{aligned}$$

and, since $w_i \geq 1/2$ for each $i \in R$,

$$\|\mathbf{x}_{R \setminus S^*}\|_2^2 \leq 2 \|\mathbf{x}_{R \setminus S^*}\|_{\mathbf{w},2}^2 \leq 2 \|\mathbf{x}_{S \setminus S^*}\|_{\mathbf{w},2}^2.$$

Therefore,

$$\begin{aligned} &\frac{4\kappa\beta s}{r-s} \|\bar{\mathbf{x}}_{R \setminus S^*}\|_2^2 \\ &\leq \frac{16\kappa\beta s}{r-s} \|\mathbf{x}_{S \setminus S^*}\|_{\mathbf{w},2}^2 + \frac{2\kappa s}{\beta(r-s)} \|\nabla_{R \setminus S^*} g(\mathbf{x})\|_2^2 \\ &\leq \frac{16\kappa\beta s}{r-s} \|\mathbf{x}_{S \setminus S^*}\|_{\mathbf{w},2}^2 + \frac{s}{2(r-s)} (g(\mathbf{x}) - f(\mathbf{x}^*)), \end{aligned}$$

where the last inequality follows from (15) since $R \setminus S^* \subseteq S$.

Plugging this back into (18), we get

$$\begin{aligned}
 & g(\mathbf{x}') - g(\mathbf{x}) \\
 & \leq -(4\kappa)^{-1} \left(\left(1 - \frac{s}{2(r-s)} \right) (g(\mathbf{x}) - f(\mathbf{x}^*)) \right. \\
 & \quad - \beta \langle \mathbf{w}\mathbf{x}, \mathbf{x}^* - \mathbf{x} \rangle - (\beta/2) \|\mathbf{x}\|_{\mathbf{w},2}^2 \\
 & \quad \left. + (\alpha/4) \|\mathbf{x}^* - \mathbf{x}\|_2^2 - \frac{16\kappa\beta s}{r-s} \|\mathbf{x}_{S \setminus S^*}\|_{\mathbf{w},2}^2 \right) \\
 & = -(4\kappa)^{-1} \left(\left(1 - \frac{s}{2(r-s)} \right) (g(\mathbf{x}) - f(\mathbf{x}^*)) \right. \\
 & \quad \underbrace{- \beta \langle \mathbf{w}\mathbf{x}_{S \cap S^*}, \mathbf{x}^* - \mathbf{x} \rangle + (\alpha/4) \|\mathbf{x}^* - \mathbf{x}\|_2^2}_{\geq -(\kappa\beta) \|\mathbf{x}_{S \cap S^*}\|_{\mathbf{w},2}^2} \\
 & \quad \left. + \beta \|\mathbf{x}_{S \setminus S^*}\|_{\mathbf{w},2}^2 - (\beta/2) \|\mathbf{x}\|_{\mathbf{w},2}^2 - \underbrace{\frac{16\kappa\beta s}{r-s} \|\mathbf{x}_{S \setminus S^*}\|_{\mathbf{w},2}^2}_{\geq -(\beta/4) \|\mathbf{x}_{S \setminus S^*}\|_{\mathbf{w},2}^2} \right) \\
 & \leq -(4\kappa)^{-1} \left(\left(1 - \frac{s}{2(r-s)} \right) (g(\mathbf{x}) - f(\mathbf{x}^*)) \right. \\
 & \quad - (\kappa\beta) \|\mathbf{x}_{S \cap S^*}\|_{\mathbf{w},2}^2 - (\beta/2) \|\mathbf{x}_{S \cap S^*}\|_{\mathbf{w},2}^2 \\
 & \quad \left. + (\beta/4) \|\mathbf{x}_{S \setminus S^*}\|_{\mathbf{w},2}^2 \right) \\
 & \leq -(4\kappa)^{-1} \left(0.5 (g(\mathbf{x}) - f(\mathbf{x}^*)) \right. \\
 & \quad \left. - (\kappa + 1)\beta \|\mathbf{x}_{S \cap S^*}\|_{\mathbf{w},2}^2 + (\beta/4) \|\mathbf{x}_{S \setminus S^*}\|_{\mathbf{w},2}^2 \right), \tag{19}
 \end{aligned}$$

where we used the fact that

$$\frac{s}{2(r-s)} \leq 1/2,$$

which holds as long as $s' \geq 4s$, and

$$\frac{16\kappa\beta s}{r-s} \leq \beta/4,$$

which holds as long as $r \geq s'/2$ and $s' \geq (128\kappa + 2)s$. In the last inequality we also used the property $\mathbf{w}/2 \leq \mathbf{w}^2$, which is by definition of \mathbf{w} .

Now, note that, because we have assumed that the first bullet of the statement doesn't hold, it has to be the case that

$$\begin{aligned}
 & (1/4) (g(\mathbf{x}) - f(\mathbf{x}^*)) \\
 & - (\kappa + 1)\beta \|\mathbf{x}_{S \cap S^*}\|_{\mathbf{w},2}^2 + (\beta/4) \|\mathbf{x}_{S \setminus S^*}\|_{\mathbf{w},2}^2 \leq 0.
 \end{aligned}$$

This immediately implies that

$$\begin{aligned}
 & (\kappa + 1)\beta \|\mathbf{x}_{S \cap S^*}\|_{\mathbf{w},2}^2 \geq (\beta/4) \|\mathbf{x}_{S \setminus S^*}\|_{\mathbf{w},2}^2 \\
 & \Rightarrow (4\kappa + 4) \|\mathbf{x}_{S \cap S^*}\|_{\mathbf{w},2}^2 \geq \|\mathbf{x}_{S \setminus S^*}\|_{\mathbf{w},2}^2 \\
 & \Rightarrow (4\kappa + 6) \|\mathbf{x}_{S \cap S^*}\|_{\mathbf{w},2}^2 \geq \|\mathbf{x}\|_{\mathbf{w},2}^2,
 \end{aligned}$$

so

$$\|\mathbf{x}_{S \cap S^*}\|_{\mathbf{w},2}^2 \geq (4\kappa + 6)^{-1} \|\mathbf{x}\|_{\mathbf{w},2}^2.$$

Similarly we also have

$$\begin{aligned}
 & (\kappa + 1)\beta \|\mathbf{x}_{S \cap S^*}\|_{\mathbf{w},2}^2 \geq (1/4) (g(\mathbf{x}) - f(\mathbf{x}^*)) \\
 & \Rightarrow (\beta/2) \|\mathbf{x}_{S \cap S^*}\|_{\mathbf{w},2}^2 \geq (8\kappa + 8)^{-1} (g(\mathbf{x}) - f(\mathbf{x}^*)).
 \end{aligned}$$

Therefore the second bullet of the statement is true, and we are done. \square

D. Low Rank Minimization

Algorithm 2 Regularized Local Search

\mathbf{A}^0 : initial rank- r' solution
 $\mathbf{W}^0 = \mathbf{Y}^0 = \mathbf{I}$: initial regularization weights
 η : step size, T : #iterations
 c : weight step size
for $t = 0, \dots, T - 1$ **do**
 $\Phi(\mathbf{A}) := (\beta/4) \left(\langle \mathbf{W}^t, \mathbf{A}\mathbf{A}^\top \rangle + \langle \mathbf{Y}^t, \mathbf{A}^\top \mathbf{A} \rangle \right)$
 $g(\mathbf{A}) := f(\mathbf{A}) + \Phi(\mathbf{A})$
 $\bar{\mathbf{A}} = H_{s'-1}(\mathbf{A}^t) - 0.5H_1(\eta \nabla g(\mathbf{A}^t))$
 $\mathbf{P} = (\mathbf{W}^t)^{1/2} \mathbf{A}^t (\mathbf{A}^t)^\top (\mathbf{W}^t)^{1/2}$
 $\mathbf{Q} = (\mathbf{Y}^t)^{1/2} (\mathbf{A}^t)^\top \mathbf{A}^t (\mathbf{Y}^t)^{1/2}$
 $\Delta = g(\mathbf{A}^t) - f(\mathbf{A}^*)$
if $g(\mathbf{A}^t) - g(\bar{\mathbf{A}}) \geq (r')^{-1} \Delta$ **then**
 Let $\mathbf{U} \mathbf{\Sigma} \mathbf{V}^\top$ be the SVD of $\bar{\mathbf{A}}$
 $\mathbf{A}^{t+1} = \operatorname{argmin}_{\mathbf{A} = \mathbf{U} \mathbf{X} \mathbf{V}^\top} g(\mathbf{A})$
 $\mathbf{W}^{t+1}, \mathbf{Y}^{t+1} = \mathbf{W}^t, \mathbf{Y}^t$
else if
 $\max \{ \operatorname{Tr}[H_r(\mathbf{P})], \operatorname{Tr}[H_r(\mathbf{Q})] \} \geq (0.4/\beta) \Delta$
then
 $\mathbf{A}^{t+1} = \mathbf{A}^t$
 $\mathbf{W}^{t+1} = (\mathbf{W}^t)^{1/2} (\mathbf{I} - r^{-1} \mathbf{I}_{\operatorname{im}(\mathbf{P})}) (\mathbf{W}^t)^{1/2}$
 $\mathbf{Y}^{t+1} = (\mathbf{Y}^t)^{1/2} (\mathbf{I} - r^{-1} \mathbf{I}_{\operatorname{im}(\mathbf{Q})}) (\mathbf{Y}^t)^{1/2}$
else
 $\mathbf{A}^{t+1} = \mathbf{A}^t$
 $\mathbf{W}^{t+1} = \mathbf{W}^t - \frac{\mathbf{W}^t \mathbf{A}^t (\mathbf{A}^t)^\top \mathbf{W}^t}{\langle \mathbf{W}^t, \mathbf{A}^t (\mathbf{A}^t)^\top \rangle}$
 $\mathbf{Y}^{t+1} = \mathbf{Y}^t - \frac{\mathbf{Y}^t (\mathbf{A}^t)^\top \mathbf{A}^t \mathbf{Y}^t}{\langle \mathbf{Y}^t, (\mathbf{A}^t)^\top \mathbf{A}^t \rangle}$
end if
end for

D.1. Preliminaries

We will use the following simple lemma about Frobenius products between low-rank projections and symmetric PSD matrices. We remind the reader that $H_r(\mathbf{A})$ is the matrix consisting of the top r components from the singular value decomposition of \mathbf{A} .

Lemma D.1. For any two symmetric PSD matrices $\mathbf{H}, \mathbf{A} \in \mathbb{R}^{n \times n}$, where $\text{rank}(\mathbf{H}) \leq r$ and $\|\mathbf{H}\|_2 \leq 1$, we have that

$$|\langle \mathbf{H}, \mathbf{A} \rangle| \leq \text{Tr}[H_r(\mathbf{A})].$$

Proof. We will use the following inequality for singular values

$$\sum_{i=1}^k \sigma_i(AB) \leq \sum_{i=1}^k \sigma_i(A)\sigma_i(B)$$

for $k = 1, \dots, n$, $A, B \in \mathbb{R}^{n \times n}$ and $\sigma_1(A) \geq \dots \geq \sigma_n(A)$ are singular values of matrix A (see page 177 in (Horn & Johnson, 1991)). Then

$$\begin{aligned} |\langle \mathbf{H}, \mathbf{A} \rangle| &= \text{Tr}[\mathbf{H} \mathbf{A}^\top] \\ &= \text{Tr}[\mathbf{H} \mathbf{A}] \\ &= \sum_{i=1}^n \sigma_i(\mathbf{H} \mathbf{A}) \\ &\leq \sum_{i=1}^n \sigma_i(\mathbf{H}) \sigma_i(\mathbf{A}) \\ &= \sum_{i=1}^r \sigma_i(\mathbf{H}) \sigma_i(\mathbf{A}) \\ &\leq \sum_{i=1}^r \sigma_i(\mathbf{A}) \\ &= \text{Tr}[H_r(\mathbf{A})]. \end{aligned}$$

□

D.2. Analysis

This section is devoted to proving Theorem 1.2, which analyzes an algorithm for low rank optimization that uses adaptive regularization.

Theorem 1.2 (Adaptive Regularization for Low Rank Optimization). Let $f \in \mathbb{R}^{m \times n} \rightarrow \mathbb{R}$ be a convex function with condition number κ and consider the low rank minimization problem

$$\min_{\text{rank}(\mathbf{A}) \leq r} f(\mathbf{A}). \quad (20)$$

For any error parameter $\varepsilon > 0$, there exists a polynomial time algorithm that returns a matrix \mathbf{A} with $\text{rank}(\mathbf{A}) \leq O\left(r \left(\kappa + \log \frac{f(\mathbf{O}) - f(\mathbf{A}^*)}{\varepsilon}\right)\right)$ and $f(\mathbf{A}) \leq f(\mathbf{A}^*) + \varepsilon$, where \mathbf{A}^* is any rank- r matrix.

Proof of Theorem 1.2. Let the smoothness and strong convexity parameters of f be β, α . We repeatedly apply Lemma D.2 $T \geq O\left(r\kappa \log \frac{f(\mathbf{A}^0) + (\beta/2)\|\mathbf{A}^0\|_F^2 - f(\mathbf{A}^*)}{\varepsilon}\right)$

times starting from solution $\mathbf{A}^0 = \mathbf{O}$ and weight matrices $\mathbf{W}^0 = \mathbf{I}, \mathbf{Y}^0 = \mathbf{I}$. Thus, we obtain solutions $\mathbf{A}^0, \dots, \mathbf{A}^T$, and weights $\mathbf{W}^0, \mathbf{W}^1, \dots, \mathbf{W}^T$ and $\mathbf{Y}^0, \mathbf{Y}^1, \dots, \mathbf{Y}^T$. We let

$$\begin{aligned} g^t(\mathbf{A}) &= f(\mathbf{A}) + (\beta/4) \left(\langle \mathbf{W}^t, \mathbf{A}^t (\mathbf{A}^t)^\top \rangle + \langle \mathbf{Y}^t, (\mathbf{A}^t)^\top \mathbf{A}^t \rangle \right) \end{aligned}$$

be the regularized function at iteration t .

We denote by T_i the total number of iterations for which item $i \in \{1, 2, 3\}$ from the statement of Lemma D.2 holds.

Consider the T_2 iterations for which item 2 from the statement of Lemma D.2 holds. Without loss of generality, \mathbf{W} is updated at least $T_2/2$ times. Letting $\mathbf{A}^* = \mathbf{U}^* \mathbf{\Sigma}^* \mathbf{V}^{*\top}$ be the singular value decomposition of \mathbf{A}^* , for each such iteration we have

$$\begin{aligned} &\text{Tr}[\mathbf{H}_{\text{im}(\mathbf{U}^*)} \mathbf{W}^{t+1} \mathbf{H}_{\text{im}(\mathbf{U}^*)}] \\ &\leq \text{Tr}[\mathbf{H}_{\text{im}(\mathbf{U}^*)} \mathbf{W}^t \mathbf{H}_{\text{im}(\mathbf{U}^*)}] - (10\kappa)^{-1}, \end{aligned}$$

and for all other types of iterations we have $\mathbf{W}^{t+1} \preceq \mathbf{W}^t$. Therefore,

$$\begin{aligned} &\text{Tr}[\mathbf{H}_{\text{im}(\mathbf{U}^*)} \mathbf{W}^T \mathbf{H}_{\text{im}(\mathbf{U}^*)}] \\ &\leq \text{Tr}[\mathbf{H}_{\text{im}(\mathbf{U}^*)} \mathbf{W}^0 \mathbf{H}_{\text{im}(\mathbf{U}^*)}] - \frac{T_2}{2} (10\kappa)^{-1}. \end{aligned}$$

However, note that by the guarantee of Lemma D.2 that $\mathbf{W}^T \succeq \mathbf{O}$, we have

$$\text{Tr}[\mathbf{H}_{\text{im}(\mathbf{U}^*)} \mathbf{W}^T \mathbf{H}_{\text{im}(\mathbf{U}^*)}] \geq 0,$$

and because $\mathbf{W}^0 = \mathbf{I}$ we also know that

$$\text{Tr}[\mathbf{H}_{\text{im}(\mathbf{U}^*)} \mathbf{W}^0 \mathbf{H}_{\text{im}(\mathbf{U}^*)}] = \text{Tr}[\mathbf{H}_{\text{im}(\mathbf{U}^*)}] \leq r.$$

This implies that $T_2 \leq 20\kappa r$.

Now, if $T_1 \geq 16r\kappa \log \frac{g^0(\mathbf{A}^0) - f(\mathbf{A}^*)}{\varepsilon}$, and since $g^t(\mathbf{A}^t)$ is non-increasing for all t , we have

$$\begin{aligned} &g^T(\mathbf{A}^T) - f(\mathbf{A}^*) \\ &\leq (1 - (16r\kappa)^{-1})^{T_1} (g^0(\mathbf{A}^0) - f(\mathbf{A}^*)) \\ &\leq \varepsilon, \end{aligned}$$

so $T_1 \leq 16r\kappa \log \frac{g^0(\mathbf{A}^0) - f(\mathbf{A}^*)}{\varepsilon}$.

Similarly, if $T_3 \geq 10r \log \frac{g^0(\mathbf{A}^0) - f(\mathbf{A}^*)}{\varepsilon}$ we have

$$\begin{aligned} &g^T(\mathbf{A}^T) - f(\mathbf{A}^*) \\ &\leq (1 - (10r)^{-1})^{T_3} (g^0(\mathbf{A}^0) - f(\mathbf{A}^*)) \\ &\leq \varepsilon, \end{aligned}$$

so $T_3 \geq 10r \log \frac{g^0(\mathbf{A}^0) - f(\mathbf{A}^*)}{\varepsilon}$.

Overall, we have that the total number of iterations is

$$T = \sum T_i \leq 36r(\kappa + 1) \log \frac{g^0(\mathbf{A}^0) - f(\mathbf{A}^*)}{\varepsilon}.$$

The only thing left is to ensure that the conditions

$$\begin{aligned} \text{Tr} [\mathbf{I} - \mathbf{W}^t] &\leq r'/2 \\ \text{Tr} [\mathbf{I} - \mathbf{Y}^t] &\leq r'/2 \end{aligned}$$

of Lemma D.2 are satisfied for all t . By the guarantees of Lemma D.2, if one of items 2, 3 holds, then

$$\text{Tr} [\mathbf{I} - \mathbf{W}^{t+1}] \leq \text{Tr} [\mathbf{I} - \mathbf{W}^t] + 1,$$

and if item 1 holds, then

$$\text{Tr} [\mathbf{I} - \mathbf{W}^{t+1}] = \text{Tr} [\mathbf{I} - \mathbf{W}^t].$$

As $\text{Tr} [\mathbf{I} - \mathbf{W}^0] = 0$, we have

$$\begin{aligned} \text{Tr} [\mathbf{I} - \mathbf{W}^T] &\leq T_2 + T_3 \\ &\leq 20\kappa r + 10r \log \frac{g^0(\mathbf{A}^0) - f(\mathbf{A}^*)}{\varepsilon} \\ &\leq r'/2, \end{aligned}$$

where the last inequality holds as long as

$$r' \geq 20r \left(2\kappa + \log \frac{g^0(\mathbf{A}^0) - f(\mathbf{A}^*)}{\varepsilon} \right).$$

□

Lemma D.2 (Low rank minimization step analysis). *Let $f : \mathbb{R}^{m \times n} \rightarrow \mathbb{R}$ be a β -smooth and α -strongly convex function with condition number $\kappa = \beta/\alpha$, and $\mathbf{W} \in \mathbb{R}^{m \times m}$, $\mathbf{Y} \in \mathbb{R}^{n \times n}$ be symmetric positive semi-definite weight matrices with spectral norm bounded by 1 and such that $\text{Tr} [\mathbf{I} - \mathbf{W}] \leq r'/2$ and $\text{Tr} [\mathbf{I} - \mathbf{Y}] \leq r'/2$ for fixed parameter $r' \geq 256r$. We define the regularized function*

$$g(\mathbf{A}) := f(\mathbf{A}) + \underbrace{(\beta/4) \left(\langle \mathbf{W}, \mathbf{A}\mathbf{A}^\top \rangle + \langle \mathbf{Y}, \mathbf{A}^\top \mathbf{A} \rangle \right)}_{\Phi(\mathbf{A})}.$$

Now, consider a rank- r' matrix $\mathbf{A} \in \mathbb{R}^{m \times n}$ with singular value decomposition

$$\mathbf{A} = \mathbf{U}\mathbf{\Lambda}\mathbf{V}^\top = \sum_{j \in S} \lambda_j \mathbf{u}_j \mathbf{v}_j^\top$$

and with the property that

$$\mathbf{\Pi}_{\text{im}(\mathbf{U})} \cdot \nabla g(\mathbf{A}) \cdot \mathbf{\Pi}_{\text{im}(\mathbf{V})} = \mathbf{O}.$$

For any rank- r solution \mathbf{A}^* where $r' \geq 256r$, there is a procedure that updates \mathbf{A} , \mathbf{W} , \mathbf{Y} , and for which exactly one of the following scenarios holds:

1. \mathbf{A} is updated to a rank- r' matrix \mathbf{A}' , and \mathbf{W} , \mathbf{Y} are not updated. We have sufficient progress in the regularized function:

$$g(\mathbf{A}') \leq g(\mathbf{A}) - (16\kappa r)^{-1} (g(\mathbf{A}) - f(\mathbf{A}^*)).$$

2. Exactly one of \mathbf{W} or \mathbf{Y} is updated (wlog \mathbf{W}) to a symmetric PSD $\mathbf{W}' \preceq \mathbf{W}$, and \mathbf{A} is not updated. We have

$$\text{Tr} [\mathbf{I} - \mathbf{W}'] \leq \text{Tr} [\mathbf{I} - \mathbf{W}] + 1$$

and

$$\begin{aligned} \text{Tr} [\mathbf{\Pi}_{\text{im}(\mathbf{U}^*)} \mathbf{W}' \mathbf{\Pi}_{\text{im}(\mathbf{U}^*)}] \\ \leq \text{Tr} [\mathbf{\Pi}_{\text{im}(\mathbf{U}^*)} \mathbf{W} \mathbf{\Pi}_{\text{im}(\mathbf{U}^*)}] - (10\kappa)^{-1}. \end{aligned}$$

Respectively, for \mathbf{Y} :

$$\begin{aligned} \text{Tr} [\mathbf{\Pi}_{\text{im}(\mathbf{V}^*)} \mathbf{Y}' \mathbf{\Pi}_{\text{im}(\mathbf{V}^*)}] \\ \leq \text{Tr} [\mathbf{\Pi}_{\text{im}(\mathbf{V}^*)} \mathbf{Y} \mathbf{\Pi}_{\text{im}(\mathbf{V}^*)}] - (10\kappa)^{-1}. \end{aligned}$$

3. Exactly one of \mathbf{W} or \mathbf{Y} is updated (wlog \mathbf{W}) to a symmetric PSD $\mathbf{W}' \preceq \mathbf{W}$, and \mathbf{A} is not updated. We have sufficient progress in the regularized function, where g' is the regularized function with the new weights:

$$g'(\mathbf{A}) \leq g(\mathbf{A}) - (10r)^{-1} (g(\mathbf{A}) - f(\mathbf{A}^*)).$$

Additionally,

$$\text{Tr} [\mathbf{I} - \mathbf{W}'] \leq \text{Tr} [\mathbf{I} - \mathbf{W}] + 1$$

Proof. We attempt to make the update $\mathbf{A} \rightarrow \mathbf{A}'$ as defined in Lemma D.3. If it makes enough progress, i.e.

$$g(\mathbf{A}') \leq g(\mathbf{A}) - (16\kappa r)^{-1} (g(\mathbf{A}) - f(\mathbf{A}^*)),$$

we are done. Otherwise, one of the items 2-5 in the statement of Lemma D.3 must hold. Let us take them one by one.

Item 2:

$$\langle \mathbf{\Pi}_{\text{im}(\mathbf{U}^*)}, \mathbf{W} \mathbf{A} \mathbf{A}^\top \mathbf{W} \rangle \geq (10\kappa)^{-1} \langle \mathbf{W}, \mathbf{A} \mathbf{A}^\top \rangle.$$

We update \mathbf{W} as

$$\mathbf{W}' = \mathbf{W} - c \cdot \mathbf{W} \mathbf{A} \mathbf{A}^\top \mathbf{W},$$

where $c = \langle \mathbf{W}, \mathbf{A} \mathbf{A}^\top \rangle^{-1}$. Note that this update preserves symmetry, and

$$\mathbf{O} \preceq \mathbf{W}' \preceq \mathbf{W}.$$

This is because

$$c \mathbf{W}^{1/2} \mathbf{A} \mathbf{A}^\top \mathbf{W}^{1/2} \preceq c \langle \mathbf{W}, \mathbf{A} \mathbf{A}^\top \rangle \cdot \mathbf{I} \preceq \mathbf{I},$$

so

$$\mathbf{W}' = \mathbf{W}^{1/2} \left(\mathbf{I} - c \mathbf{W}^{1/2} \mathbf{A} \mathbf{A}^\top \mathbf{W}^{1/2} \right) \mathbf{W}^{1/2} \succeq \mathbf{O}$$

and

$$\mathbf{W}' = \mathbf{W} - c \mathbf{W} \mathbf{A} \mathbf{A}^\top \mathbf{W} \preceq \mathbf{W}.$$

Now, note that

$$\begin{aligned} \text{Tr} [\mathbf{I} - \mathbf{W}'] &= \text{Tr} [\mathbf{I} - \mathbf{W}] + c \langle \mathbf{W}^2, \mathbf{A} \mathbf{A}^\top \rangle \\ &\leq \text{Tr} [\mathbf{I} - \mathbf{W}] + c \langle \mathbf{W}, \mathbf{A} \mathbf{A}^\top \rangle \\ &= \text{Tr} [\mathbf{I} - \mathbf{W}] + 1, \end{aligned}$$

where we used the fact that $\mathbf{W}^2 \preceq \mathbf{W}$, and (letting $\mathbf{\Pi}^* = \mathbf{\Pi}_{\text{im}(U^*)}$ for convenience),

$$\begin{aligned} \text{Tr} [\mathbf{\Pi}^* \mathbf{W}' \mathbf{\Pi}^*] &= \text{Tr} [\mathbf{\Pi}^* \mathbf{W} \mathbf{\Pi}^*] - c \langle \mathbf{\Pi}^*, \mathbf{W} \mathbf{A} \mathbf{A}^\top \mathbf{W} \rangle \\ &\leq \text{Tr} [\mathbf{\Pi}^* \mathbf{W} \mathbf{\Pi}^*] - c / (10\kappa) \langle \mathbf{W}, \mathbf{A} \mathbf{A}^\top \rangle \\ &= \text{Tr} [\mathbf{\Pi}^* \mathbf{W} \mathbf{\Pi}^*] - (10\kappa)^{-1}, \end{aligned} \quad (21)$$

Item 3:

$$\langle \mathbf{\Pi}_{\text{im}(V^*)}, \mathbf{Y} \mathbf{A}^\top \mathbf{A} \mathbf{Y} \rangle \geq (10\kappa)^{-1} \langle \mathbf{Y}, \mathbf{A}^\top \mathbf{A} \rangle.$$

This is entirely analogous to the previous case.

Item 4:

$$(\beta/4) \text{Tr} \left[H_r \left(\mathbf{A}^\top \mathbf{W} \mathbf{A} \right) \right] \geq 10^{-1} (g(\mathbf{A}) - f(\mathbf{A}^*)). \quad (22)$$

After considering the eigendecomposition

$$\mathbf{W}^{1/2} \mathbf{A} \mathbf{A}^\top \mathbf{W}^{1/2} = \sum_{i \in [r']} \bar{\lambda}_i \bar{\mathbf{v}}_i \bar{\mathbf{v}}_i^\top$$

with $\bar{\lambda}_1 \geq \bar{\lambda}_2 \geq \dots \geq \bar{\lambda}_{r'} \geq 0$, (22) can be re-phrased as

$$(\beta/4) \sum_{i \in [r]} \bar{\lambda}_i > (1/10) (g(\mathbf{A}) - f(\mathbf{A}^*)).$$

We update \mathbf{W} as

$$\mathbf{W}' = \mathbf{W}^{1/2} \left(\mathbf{I} - r^{-1} \sum_{i \in [r]} \bar{\mathbf{v}}_i \bar{\mathbf{v}}_i^\top \right) \mathbf{W}^{1/2}$$

and let g' be the new regularized objective. First of all, note that this operation preserves symmetry, and that $\mathbf{O} \preceq \mathbf{W}' \preceq \mathbf{I}$, since $\sum_{i \in [r]} \bar{\mathbf{v}}_i \bar{\mathbf{v}}_i^\top \preceq \mathbf{I}$. Additionally,

$$\begin{aligned} \text{Tr} [\mathbf{I} - \mathbf{W}'] &= \text{Tr} [\mathbf{I} - \mathbf{W}] + r^{-1} \sum_{i \in [r]} \bar{\mathbf{v}}_i^\top \mathbf{W} \bar{\mathbf{v}}_i \\ &\leq \text{Tr} [\mathbf{W}] + 1 \end{aligned}$$

and

$$\begin{aligned} g'(\mathbf{A}) - g(\mathbf{A}) &= (\beta/4) \langle \mathbf{W}', \mathbf{A} \mathbf{A}^\top \rangle - (\beta/4) \langle \mathbf{W}, \mathbf{A} \mathbf{A}^\top \rangle \\ &= -(\beta/(4r)) \left\langle \mathbf{W}^{1/2} \left(\sum_{i \in [r]} \bar{\mathbf{v}}_i \bar{\mathbf{v}}_i^\top \right) \mathbf{W}^{1/2}, \mathbf{A} \mathbf{A}^\top \right\rangle \\ &= -(\beta/(4r)) \sum_{i \in [r]} \bar{\lambda}_i \\ &\leq -(10r)^{-1} (g(\mathbf{A}) - f(\mathbf{A}^*)), \end{aligned}$$

Item 5:

$$(\beta/4) \text{Tr} \left[H_r \left(\mathbf{A} \mathbf{Y} \mathbf{A}^\top \right) \right] \geq 10^{-1} (g(\mathbf{A}) - f(\mathbf{A}^*)).$$

This is entirely analogous to the previous case. \square

Lemma D.3. Let $f : \mathbb{R}^{m \times n} \rightarrow \mathbb{R}$ be a β -smooth and α -strongly convex function with condition number $\kappa = \beta/\alpha$, and $\mathbf{W} \in \mathbb{R}^{m \times m}$, $\mathbf{Y} \in \mathbb{R}^{n \times n}$ be symmetric positive semi-definite weight matrices with spectral norm bounded by 1 and such that $\text{Tr} [\mathbf{I} - \mathbf{W}]$, $\text{Tr} [\mathbf{I} - \mathbf{Y}] \leq r'/2$ for some parameter $r' \geq 0$. We define the regularized function

$$g(\mathbf{A}) := f(\mathbf{A}) + (\beta/4) \underbrace{\left(\langle \mathbf{W}, \mathbf{A} \mathbf{A}^\top \rangle + \langle \mathbf{Y}, \mathbf{A}^\top \mathbf{A} \rangle \right)}_{\Phi(\mathbf{A})}.$$

Now, consider a rank- r' matrix $\mathbf{A} \in \mathbb{R}^{m \times n}$ with singular value decomposition

$$\mathbf{A} = \mathbf{U} \mathbf{\Lambda} \mathbf{V}^\top = \sum_{j \in S} \lambda_j \mathbf{u}_j \mathbf{v}_j^\top$$

and with the property that

$$\mathbf{\Pi}_{\text{im}(U)} \cdot \nabla g(\mathbf{A}) \cdot \mathbf{\Pi}_{\text{im}(V)} = \mathbf{O}.$$

We define an updated solution

$$\mathbf{A}' = \mathbf{A} - \eta \cdot H_1(\nabla g(\mathbf{A})) - \lambda_j \mathbf{u}_j \mathbf{v}_j^\top,$$

where $\eta = (2\beta)^{-1}$, $H_1(\cdot)$ returns the top singular component, and $j \in S$ is picked to minimize λ_j .

Then, for any rank- r solution \mathbf{A}^* , where $r' \geq 256r$, and its singular value decomposition $\mathbf{A}^* = \mathbf{U}^* \mathbf{\Lambda}^* \mathbf{V}^{*\top}$, at least one of the following conditions holds:

1. We have sufficient progress in the regularized function:

$$g(\mathbf{A}') \leq g(\mathbf{A}) - (16\kappa r)^{-1} (g(\mathbf{A}) - f(\mathbf{A}^*)).$$

2. $\mathbf{W} \mathbf{A} \mathbf{A}^\top \mathbf{W}$ is significantly correlated to \mathbf{U}^* :

$$\langle \mathbf{\Pi}_{\text{im}(U^*)}, \mathbf{W} \mathbf{A} \mathbf{A}^\top \mathbf{W} \rangle \geq (10\kappa)^{-1} \langle \mathbf{W}, \mathbf{A} \mathbf{A}^\top \rangle.$$

3. $\mathbf{Y} \mathbf{A}^\top \mathbf{A} \mathbf{Y}$ is significantly correlated to \mathbf{V}^* :

$$\langle \mathbf{I}_{\text{im}(\mathbf{V}^*)}, \mathbf{Y} \mathbf{A}^\top \mathbf{A} \mathbf{Y} \rangle \geq (10\kappa)^{-1} \langle \mathbf{Y}, \mathbf{A}^\top \mathbf{A} \rangle.$$

4. The spectrum of $\mathbf{A}^\top \mathbf{W} \mathbf{A}$ is highly concentrated and responsible for a constant fraction of the error:

$$(\beta/4) \text{Tr} \left[H_r \left(\mathbf{A}^\top \mathbf{W} \mathbf{A} \right) \right] \geq 10^{-1} (g(\mathbf{A}) - f(\mathbf{A}^*)).$$

and

5. The spectrum of $\mathbf{A} \mathbf{Y} \mathbf{A}^\top$ is highly concentrated and responsible for a constant fraction of the error:

$$(\beta/4) \text{Tr} \left[H_r \left(\mathbf{A} \mathbf{Y} \mathbf{A}^\top \right) \right] \geq 10^{-1} (g(\mathbf{A}) - f(\mathbf{A}^*)).$$

Proof. Note that g is a 2β -smooth function. This follows because

$$\nabla g(\mathbf{A}) = \nabla f(\mathbf{A}) + (\beta/2) (\mathbf{W} \mathbf{A} + \mathbf{A} \mathbf{Y}),$$

and so for any two matrices \mathbf{A}, \mathbf{A}' ,

$$\begin{aligned} & \|\nabla g(\mathbf{A}') - \nabla g(\mathbf{A})\|_F \\ & \leq \|\nabla f(\mathbf{A}') - \nabla f(\mathbf{A})\|_F \\ & \quad + (\beta/2) \|\mathbf{W}(\mathbf{A}' - \mathbf{A})\|_F + (\beta/2) \|(\mathbf{A}' - \mathbf{A})\mathbf{Y}\|_F \\ & \leq 2\beta \|\mathbf{A}' - \mathbf{A}\|_F, \end{aligned}$$

which is known to imply 2β -smoothness of g . Here we used the triangle inequality and the fact that $\mathbf{W}, \mathbf{Y} \preceq \mathbf{I}$. Therefore, we have

$$\begin{aligned} & g(\mathbf{A}') - g(\mathbf{A}) \\ & \leq \langle \nabla g(\mathbf{A}), \mathbf{A}' - \mathbf{A} \rangle + \|\nabla g(\mathbf{A}') - \nabla g(\mathbf{A})\|_F \|\mathbf{A}' - \mathbf{A}\|_F \\ & \leq \langle \nabla g(\mathbf{A}), \mathbf{A}' - \mathbf{A} \rangle + \beta \|\mathbf{A}' - \mathbf{A}\|_F^2 \\ & \leq -\eta \|\nabla g(\mathbf{A})\|_2^2 + 2\beta\eta^2 \|\nabla g(\mathbf{A})\|_2^2 + 2\beta\lambda_j^2 \\ & = -(8\beta)^{-1} \|\nabla g(\mathbf{A})\|_2^2 + 2\beta\lambda_j^2, \end{aligned} \tag{23}$$

where in the second inequality we used the facts that

$$\begin{aligned} & \langle \nabla g(\mathbf{A}), -\lambda_j \mathbf{u}_j \mathbf{v}_j^\top \rangle \\ & = \langle \mathbf{I}_{\text{im}(\mathbf{U})} \nabla g(\mathbf{A}) \mathbf{I}_{\text{im}(\mathbf{V})}, -\lambda_j \mathbf{u}_j \mathbf{v}_j^\top \rangle \\ & = 0 \end{aligned}$$

and that, for any two matrices \mathbf{B}, \mathbf{C} ,

$$\|\mathbf{B} + \mathbf{C}\|_F^2 \leq 2\|\mathbf{B}\|_F^2 + 2\|\mathbf{C}\|_F^2.$$

The last equality follows by our choice of η . In order to lower bound $\|\nabla g(\mathbf{A})\|_2^2$, we use the strong convexity of f

as follows:

$$\begin{aligned} & f(\mathbf{A}^*) - f(\mathbf{A}) \\ & \geq \langle \nabla f(\mathbf{A}), \mathbf{A}^* - \mathbf{A} \rangle + (\alpha/2) \|\mathbf{A}^* - \mathbf{A}\|_F^2 \\ & = \langle \nabla g(\mathbf{A}), \mathbf{A}^* - \mathbf{A} \rangle \\ & \quad - \langle \nabla \Phi(\mathbf{A}), \mathbf{A}^* - \mathbf{A} \rangle + (\alpha/2) \|\mathbf{A}^* - \mathbf{A}\|_F^2 \\ & = \underbrace{\langle \nabla g(\mathbf{A}), \mathbf{A}^* - \mathbf{A} \rangle + (\alpha/4) \|\mathbf{A}^* - \mathbf{A}\|_F^2}_P \\ & \quad - \langle \nabla \Phi(\mathbf{A}), \mathbf{A}^* - \mathbf{A} \rangle + (\alpha/4) \|\mathbf{A}^* - \mathbf{A}\|_F^2. \end{aligned} \tag{24}$$

Bounding P . We let $\mathbf{I}_{\text{im}(\mathbf{U})}, \mathbf{I}_{\text{im}(\mathbf{V})}$ be the orthogonal projections onto the images of \mathbf{U} and \mathbf{V} respectively, so we can write

$$\begin{aligned} & \mathbf{A}^* - \mathbf{A} \\ & = \mathbf{I}_{\text{im}(\mathbf{U})} (\mathbf{A}^* - \mathbf{A}) \mathbf{I}_{\text{im}(\mathbf{V})} \\ & \quad + (\mathbf{I} - \mathbf{I}_{\text{im}(\mathbf{U})}) (\mathbf{A}^* - \mathbf{A}) \mathbf{I}_{\text{im}(\mathbf{V})} \\ & \quad + (\mathbf{A}^* - \mathbf{A}) (\mathbf{I} - \mathbf{I}_{\text{im}(\mathbf{V})}) \\ & = \mathbf{I}_{\text{im}(\mathbf{U})} (\mathbf{A}^* - \mathbf{A}) \mathbf{I}_{\text{im}(\mathbf{V})} \\ & \quad + (\mathbf{I} - \mathbf{I}_{\text{im}(\mathbf{U})}) \mathbf{A}^* \mathbf{I}_{\text{im}(\mathbf{V})} \\ & \quad + \mathbf{A}^* (\mathbf{I} - \mathbf{I}_{\text{im}(\mathbf{V})}). \end{aligned}$$

Now, note that

$$\begin{aligned} & \langle \nabla g(\mathbf{A}), \mathbf{A}^* - \mathbf{A} \rangle \\ & = \langle \nabla g(\mathbf{A}), (\mathbf{I} - \mathbf{I}_{\text{im}(\mathbf{U})}) \mathbf{A}^* \mathbf{I}_{\text{im}(\mathbf{V})} \rangle \\ & \quad + \langle \nabla g(\mathbf{A}), \mathbf{A}^* (\mathbf{I} - \mathbf{I}_{\text{im}(\mathbf{V})}) \rangle, \end{aligned}$$

where we used the fact that

$$\begin{aligned} & \langle \nabla g(\mathbf{A}), \mathbf{I}_{\text{im}(\mathbf{U})} (\mathbf{A}^* - \mathbf{A}) \mathbf{I}_{\text{im}(\mathbf{V})} \rangle \\ & = \langle \mathbf{I}_{\text{im}(\mathbf{U})} \nabla g(\mathbf{A}) \mathbf{I}_{\text{im}(\mathbf{V})}, \mathbf{A}^* - \mathbf{A} \rangle \\ & = 0, \end{aligned}$$

and

$$\begin{aligned} & (\alpha/4) \|\mathbf{A}^* - \mathbf{A}\|_F^2 \\ & \geq (\alpha/4) \|(\mathbf{I} - \mathbf{I}_{\text{im}(\mathbf{U})}) \mathbf{A}^* \mathbf{I}_{\text{im}(\mathbf{V})}\|_F^2 \\ & \quad + (\alpha/4) \|\mathbf{A}^* (\mathbf{I} - \mathbf{I}_{\text{im}(\mathbf{V})})\|_F^2. \end{aligned}$$

Additionally, note that for any rank- r matrix \mathbf{B} , we have

$$\begin{aligned} & \langle \nabla g(\mathbf{A}), \mathbf{B} \rangle + (\alpha/4) \|\mathbf{B}\|_F^2 \\ & \geq -\alpha^{-1} \|H_r(\nabla g(\mathbf{A}))\|_F^2 \\ & \geq -\alpha^{-1} r \|\nabla g(\mathbf{A})\|_2^2, \end{aligned}$$

a proof of which can be found e.g. in Lemma A.6 of (Axiotis & Sviridenko, 2021a). Applying this inequality with

$$\mathbf{B} = (\mathbf{I} - \mathbf{I}_{\text{im}(\mathbf{U})}) \mathbf{A}^* \mathbf{I}_{\text{im}(\mathbf{V})}$$

and

$$\mathbf{B} = \mathbf{A}^* (\mathbf{I} - \mathbf{\Pi}_{\text{im}(\mathbf{V})})$$

and summing them up, we obtain

$$\begin{aligned} P &= \langle \nabla g(\mathbf{A}), \mathbf{A}^* - \mathbf{A} \rangle + (\alpha/4) \|\mathbf{A}^* - \mathbf{A}\|_F^2 \\ &\geq -2\alpha^{-1}r \|\nabla g(\mathbf{A})\|_2^2. \end{aligned}$$

Plugging this into (24) and re-arranging, we get

$$\begin{aligned} &\|\nabla g(\mathbf{A})\|_2^2 \\ &\geq \alpha/(2r) \left(f(\mathbf{A}) - f(\mathbf{A}^*) \right. \\ &\quad \left. - \langle \nabla \Phi(\mathbf{A}), \mathbf{A}^* - \mathbf{A} \rangle + (\alpha/4) \|\mathbf{A}^* - \mathbf{A}\|_F^2 \right) \\ &= \alpha/(2r) \left(g(\mathbf{A}) - f(\mathbf{A}^*) - \Phi(\mathbf{A}) \right. \\ &\quad \left. - \underbrace{\langle \nabla \Phi(\mathbf{A}), \mathbf{A}^* - \mathbf{A} \rangle + (\alpha/4) \|\mathbf{A}^* - \mathbf{A}\|_F^2}_{Q} \right). \end{aligned} \quad (25)$$

Bounding Q . We know that

$$-\langle \nabla \Phi(\mathbf{A}), \mathbf{A}^* - \mathbf{A} \rangle = -(\beta/2) \langle \mathbf{W}\mathbf{A} + \mathbf{A}\mathbf{Y}, \mathbf{A}^* - \mathbf{A} \rangle.$$

If we let

$$\mathbf{A}^* = \mathbf{U}^* \mathbf{\Lambda}^* \mathbf{V}^{*\top}$$

be the SVD of \mathbf{A}^* and $\mathbf{\Pi}_{\text{im}(\mathbf{U}^*)}$, $\mathbf{\Pi}_{\text{im}(\mathbf{V}^*)}$ be the orthogonal projections onto the images of \mathbf{U}^* and \mathbf{V}^* respectively, then we have

$$\begin{aligned} &-(\beta/2) \langle \mathbf{W}\mathbf{A}, \mathbf{A}^* - \mathbf{A} \rangle \\ &= -(\beta/2) \langle \mathbf{W}\mathbf{A}, \mathbf{\Pi}_{\text{im}(\mathbf{U}^*)}(\mathbf{A}^* - \mathbf{A}) \mathbf{\Pi}_{\text{im}(\mathbf{V}^*)} \rangle \\ &+ (\beta/2) \langle \mathbf{W}, \mathbf{A}\mathbf{A}^\top \rangle - (\beta/2) \langle \mathbf{W}\mathbf{A}, \mathbf{\Pi}_{\text{im}(\mathbf{U}^*)} \mathbf{A} \mathbf{\Pi}_{\text{im}(\mathbf{V}^*)} \rangle. \end{aligned}$$

Looking at the first term of this, we have

$$\begin{aligned} &-(\beta/2) \langle \mathbf{W}\mathbf{A}, \mathbf{\Pi}_{\text{im}(\mathbf{U}^*)}(\mathbf{A}^* - \mathbf{A}) \mathbf{\Pi}_{\text{im}(\mathbf{V}^*)} \rangle \\ &+ (\alpha/8) \|\mathbf{A}^* - \mathbf{A}\|_F^2 \\ &= -(\beta/2) \langle \mathbf{\Pi}_{\text{im}(\mathbf{U}^*)} \mathbf{W} \mathbf{A} \mathbf{\Pi}_{\text{im}(\mathbf{V}^*)}, \mathbf{A}^* - \mathbf{A} \rangle \\ &+ (\alpha/8) \|\mathbf{A}^* - \mathbf{A}\|_F^2 \\ &\geq -\beta^2/(2\alpha) \left\| \mathbf{\Pi}_{\text{im}(\mathbf{U}^*)} \mathbf{W} \mathbf{A} \mathbf{\Pi}_{\text{im}(\mathbf{V}^*)} \right\|_F^2. \end{aligned}$$

Similarly for the terms containing \mathbf{Y} , we get

$$\begin{aligned} &-(\beta/2) \langle \mathbf{A}\mathbf{Y}, \mathbf{A}^* - \mathbf{A} \rangle \\ &= -(\beta/2) \langle \mathbf{A}\mathbf{Y}, \mathbf{\Pi}_{\text{im}(\mathbf{U}^*)}(\mathbf{A}^* - \mathbf{A}) \mathbf{\Pi}_{\text{im}(\mathbf{V}^*)} \rangle \\ &+ (\beta/2) \langle \mathbf{Y}, \mathbf{A}^\top \mathbf{A} \rangle - (\beta/2) \langle \mathbf{A}\mathbf{Y}, \mathbf{\Pi}_{\text{im}(\mathbf{U}^*)} \mathbf{A} \mathbf{\Pi}_{\text{im}(\mathbf{V}^*)} \rangle. \end{aligned}$$

and

$$\begin{aligned} &-(\beta/2) \langle \mathbf{A}\mathbf{Y}, \mathbf{\Pi}_{\text{im}(\mathbf{U}^*)}(\mathbf{A}^* - \mathbf{A}) \mathbf{\Pi}_{\text{im}(\mathbf{V}^*)} \rangle \\ &+ (\alpha/8) \|\mathbf{A}^* - \mathbf{A}\|_F^2 \\ &\geq -\beta^2/(2\alpha) \left\| \mathbf{\Pi}_{\text{im}(\mathbf{U}^*)} \mathbf{A}\mathbf{Y} \mathbf{\Pi}_{\text{im}(\mathbf{V}^*)} \right\|_F^2. \end{aligned}$$

In summary, we have

$$\begin{aligned} Q &= -\langle \nabla \Phi(\mathbf{A}), \mathbf{A}^* - \mathbf{A} \rangle + (\alpha/4) \|\mathbf{A}^* - \mathbf{A}\|_F^2 \\ &\geq -\beta^2/(2\alpha) \left\| \mathbf{\Pi}_{\text{im}(\mathbf{U}^*)} \mathbf{W} \mathbf{A} \mathbf{\Pi}_{\text{im}(\mathbf{V}^*)} \right\|_F^2 \\ &\quad - \beta^2/(2\alpha) \left\| \mathbf{\Pi}_{\text{im}(\mathbf{U}^*)} \mathbf{A}\mathbf{Y} \mathbf{\Pi}_{\text{im}(\mathbf{V}^*)} \right\|_F^2 \\ &\quad + (\beta/2) \langle \mathbf{W}, \mathbf{A}\mathbf{A}^\top \rangle + (\beta/2) \langle \mathbf{Y}, \mathbf{A}^\top \mathbf{A} \rangle \\ &\quad - (\beta/2) \langle \mathbf{W}\mathbf{A}, \mathbf{\Pi}_{\text{im}(\mathbf{U}^*)} \mathbf{A} \mathbf{\Pi}_{\text{im}(\mathbf{V}^*)} \rangle \\ &\quad - (\beta/2) \langle \mathbf{A}\mathbf{Y}, \mathbf{\Pi}_{\text{im}(\mathbf{U}^*)} \mathbf{A} \mathbf{\Pi}_{\text{im}(\mathbf{V}^*)} \rangle. \end{aligned} \quad (26)$$

Now, let us assume that all items 2-5 from the lemma statement are false. For the first term of (26), we have

$$\begin{aligned} &-\beta^2/(2\alpha) \left\| \mathbf{\Pi}_{\text{im}(\mathbf{U}^*)} \mathbf{W} \mathbf{A} \mathbf{\Pi}_{\text{im}(\mathbf{V}^*)} \right\|_F^2 \\ &\geq -\beta^2/(2\alpha) \left\| \mathbf{\Pi}_{\text{im}(\mathbf{U}^*)} \mathbf{W} \mathbf{A} \right\|_F^2 \\ &= -\beta^2/(2\alpha) \langle \mathbf{\Pi}_{\text{im}(\mathbf{U}^*)}, \mathbf{W} \mathbf{A} \mathbf{A}^\top \mathbf{W} \rangle \\ &\geq -(\beta/20) \langle \mathbf{W}, \mathbf{A}\mathbf{A}^\top \rangle, \end{aligned}$$

where we used item 2 from the lemma statement, and similarly for the second term of (26),

$$\begin{aligned} &-\beta^2/(2\alpha) \left\| \mathbf{\Pi}_{\text{im}(\mathbf{U}^*)} \mathbf{A}\mathbf{Y} \mathbf{\Pi}_{\text{im}(\mathbf{V}^*)} \right\|_F^2 \\ &\geq -(\beta/20) \langle \mathbf{Y}, \mathbf{A}^\top \mathbf{A} \rangle. \end{aligned}$$

Now we look at the second to last term of (26), i.e.

$$\begin{aligned} &-(\beta/2) \langle \mathbf{W}\mathbf{A}, \mathbf{\Pi}_{\text{im}(\mathbf{U}^*)} \mathbf{A} \mathbf{\Pi}_{\text{im}(\mathbf{V}^*)} \rangle \\ &= -(\beta/2) \langle \mathbf{W} \mathbf{A} \mathbf{\Pi}_{\text{im}(\mathbf{V}^*)} \mathbf{A}^\top, \mathbf{\Pi}_{\text{im}(\mathbf{U}^*)} \rangle. \end{aligned}$$

Now, we use the matrix Holder inequality

$$\begin{aligned} &-(\beta/2) \langle \mathbf{W} \mathbf{A} \mathbf{\Pi}_{\text{im}(\mathbf{V}^*)} \mathbf{A}^\top, \mathbf{\Pi}_{\text{im}(\mathbf{U}^*)} \rangle \\ &\geq -(\beta/2) \left\| \mathbf{W} \mathbf{A} \mathbf{\Pi}_{\text{im}(\mathbf{V}^*)} \mathbf{A}^\top \right\|_* \left\| \mathbf{\Pi}_{\text{im}(\mathbf{U}^*)} \right\|_2 \\ &\geq -(\beta/2) \left\| \mathbf{W} \mathbf{A} \mathbf{\Pi}_{\text{im}(\mathbf{V}^*)} \mathbf{A}^\top \right\|_*, \end{aligned}$$

which can be proved by applying von Neumann's trace inequality and then the classical Holder inequality. Now, note that the matrix $\mathbf{W} \mathbf{A} \mathbf{\Pi}_{\text{im}(\mathbf{V}^*)} \mathbf{A}^\top$ is similar to $\mathbf{W}^{1/2} \mathbf{A} \mathbf{\Pi}_{\text{im}(\mathbf{V}^*)} \mathbf{A}^\top \mathbf{W}^{1/2}$, and so they have the same eigenvalues. Furthermore, the latter is a symmetric PSD matrix, and so the former has real positive eigenvalues as well. This means that its singular values are the same as its eigenvalues, and as a result the nuclear norm is equal to the trace, i.e.

$$\begin{aligned} &-(\beta/2) \left\| \mathbf{W} \mathbf{A} \mathbf{\Pi}_{\text{im}(\mathbf{V}^*)} \mathbf{A}^\top \right\|_* \\ &= -(\beta/2) \text{Tr} \left(\mathbf{W} \mathbf{A} \mathbf{\Pi}_{\text{im}(\mathbf{V}^*)} \mathbf{A}^\top \right) \\ &= -(\beta/2) \langle \mathbf{\Pi}_{\text{im}(\mathbf{V}^*)}, \mathbf{A}^\top \mathbf{W} \mathbf{A} \rangle \\ &\geq -(\beta/2) \text{Tr} \left[H_r \left(\mathbf{A}^\top \mathbf{W} \mathbf{A} \right) \right] \\ &\geq -(1/5) (g(\mathbf{A}) - f(\mathbf{A}^*)). \end{aligned}$$

where we also used Lemma D.1 and item 4 from the lemma statement. So we derived that

$$\begin{aligned} & -(\beta/2)\langle \mathbf{W}\mathbf{A}, \mathbf{I}\mathbf{I}_{\text{im}(\mathbf{U}^*)}\mathbf{A}\mathbf{I}\mathbf{I}_{\text{im}(\mathbf{V}^*)}\rangle \\ & \geq -(1/5)(g(\mathbf{A}) - f(\mathbf{A}^*)), \end{aligned}$$

and similarly for the last term of (26),

$$\begin{aligned} & -(\beta/2)\langle \mathbf{A}\mathbf{Y}, \mathbf{I}\mathbf{I}_{\text{im}(\mathbf{U}^*)}\mathbf{A}\mathbf{I}\mathbf{I}_{\text{im}(\mathbf{V}^*)}\rangle \\ & \geq -(1/5)(g(\mathbf{A}) - f(\mathbf{A}^*)). \end{aligned}$$

Plugging the four inequalities that we derived back into (26), we get

$$\begin{aligned} Q & \geq (\beta/2 - \beta/20)\langle \mathbf{W}, \mathbf{A}\mathbf{A}^\top \rangle \\ & + (\beta/2 - \beta/20)\langle \mathbf{Y}, \mathbf{A}^\top \mathbf{A} \rangle \\ & - (2/5)(g(\mathbf{A}) - f(\mathbf{A}^*)) \\ & = (9/5)\Phi(\mathbf{A}) - (2/5)(g(\mathbf{A}) - f(\mathbf{A}^*)) \\ & > (3/2)\Phi(\mathbf{A}) - (2/5)(g(\mathbf{A}) - f(\mathbf{A}^*)). \end{aligned}$$

Finally, combining this with the smoothness inequality (23) and the lower bound on $\|\nabla g(\mathbf{A})\|_2^2$ (25), we derive

$$\begin{aligned} g(\mathbf{A}') - g(\mathbf{A}) & \leq -(16\kappa r)^{-1}(g(\mathbf{A}) - f(\mathbf{A}^*) + (1/2)\Phi(\mathbf{A})) + 2\beta\lambda_j^2 \\ & = -(16\kappa r)^{-1}(g(\mathbf{A}) - f(\mathbf{A}^*)) - (32\kappa r)^{-1}\Phi(\mathbf{A}) + 2\beta\lambda_j^2. \end{aligned}$$

What remains is the bound the sum of the last two terms. We remind the reader that $\mathbf{A} = \mathbf{U}\mathbf{\Lambda}\mathbf{V}^\top$. Now, letting \mathbf{z} equal to the vectorized diagonal of $\mathbf{U}^\top \mathbf{W}\mathbf{U}$ and $\boldsymbol{\lambda}$ to the vectorized diagonal of $\mathbf{\Lambda}$, note that

$$\|\boldsymbol{\lambda}\|_z^2 = \langle \mathbf{\Lambda}^2, \mathbf{U}^\top \mathbf{W}\mathbf{U} \rangle = \langle \mathbf{W}, \mathbf{A}\mathbf{A}^\top \rangle,$$

using which we derive

$$\begin{aligned} \lambda_j^2 & = \min_{j \in S} \lambda_j^2 \leq \frac{\|\boldsymbol{\lambda}\|_z^2}{\|\mathbf{z}\|_1} \\ & = \frac{\langle \mathbf{W}, \mathbf{A}\mathbf{A}^\top \rangle}{\text{Tr}[\mathbf{U}^\top \mathbf{W}\mathbf{U}]} \\ & = \frac{\langle \mathbf{W}, \mathbf{A}\mathbf{A}^\top \rangle}{\text{Tr}[\mathbf{U}^\top \mathbf{U}] - \text{Tr}[\mathbf{U}^\top (\mathbf{I} - \mathbf{W})\mathbf{U}]} \\ & \leq \frac{\langle \mathbf{W}, \mathbf{A}\mathbf{A}^\top \rangle}{r' - \text{Tr}[\mathbf{I} - \mathbf{W}]} \\ & \leq \frac{\langle \mathbf{W}, \mathbf{A}\mathbf{A}^\top \rangle}{r'/2} \\ & \leq \frac{\langle \mathbf{W}, \mathbf{A}\mathbf{A}^\top \rangle}{128r\kappa}, \end{aligned}$$

where we used the fact that

$$\begin{aligned} & \text{Tr}[\mathbf{U}^\top (\mathbf{I} - \mathbf{W})\mathbf{U}] \\ & = \text{Tr}[(\mathbf{I} - \mathbf{W})^{1/2} \mathbf{U}\mathbf{U}^\top (\mathbf{I} - \mathbf{W})^{1/2}] \\ & \leq \text{Tr}[\mathbf{I} - \mathbf{W}], \end{aligned}$$

because the columns of \mathbf{U} are orthonormal. We also used the property that $\text{Tr}[\mathbf{I} - \mathbf{W}] \leq r'/2$ and the fact that $r' \geq 256r\kappa$ by the lemma statement.

Similarly, we derive that

$$\lambda_j^2 \leq \frac{\langle \mathbf{Y}, \mathbf{A}^\top \mathbf{A} \rangle}{128r\kappa},$$

and, adding these two inequalities, we have

$$2\beta\lambda_j^2 \leq (32r\kappa)^{-1}\Phi(\mathbf{A}),$$

finally concluding that

$$g(\mathbf{A}') - g(\mathbf{A}) \leq -(16\kappa r)^{-1}(g(\mathbf{A}) - f(\mathbf{A}^*)). \quad \square$$

E. Lower Bounds

Lemma E.1 (IHT lower bound). *Let $f(\mathbf{x}) := (1/2)\|\mathbf{A}\mathbf{x} - \mathbf{b}\|_2^2$. For any $\kappa, s \geq 1$, $s' \leq 0.6s\kappa^2$, there exists a (diagonal) matrix $\mathbf{A} \in \mathbb{R}^{n \times n}$ and a vector $\mathbf{b} \in \mathbb{R}^n$ where $n = s(\kappa^2 + \kappa + 1)$, f is 1-strongly convex and κ -smooth, as well as an s -sparse solution \mathbf{x}^* and an s' -sparse solution \mathbf{x} , such that*

$$f(\mathbf{x}) \geq f(\mathbf{x}^*) + 0.1s\kappa^2$$

but

$$\mathbf{x} = H_{s'}(\mathbf{x} - \beta^{-1}\nabla f(\mathbf{x})),$$

i.e. \mathbf{x} is a fixpoint for IHT.

Proof. We use the same example as in (Axiotis & Sviridenko, 2021b), Section 5.2: \mathbf{A} is diagonal with

$$\mathbf{A}_{ii} = \begin{cases} 1 & \text{if } i \in I_1 \\ \sqrt{\kappa} & \text{if } i \in I_2 \\ 1 & \text{if } i \in I_3, \end{cases}$$

where $I_1 = [s]$, $I_2 = [s+1, s(\kappa+1)]$, $I_3 = [s(\kappa+1) + 1, s(\kappa^2 + \kappa + 1)]$, and \mathbf{b} is defined as

$$\mathbf{b}_i = \begin{cases} \kappa\sqrt{1-4\delta} & \text{if } i \in I_1 \\ \sqrt{\kappa}\sqrt{1-2\delta} & \text{if } i \in I_2 \\ 1 & \text{if } i \in I_3, \end{cases}$$

for some sufficiently small $\delta > 0$ used for tie-breaking. We define

$$\mathbf{x}_i^* = \begin{cases} \kappa\sqrt{1-4\delta} & \text{if } i \in I_1 \\ 0 & \text{otherwise} \end{cases}$$

and, for some arbitrary s' -sized $S \subseteq I_3$

$$x_i = \begin{cases} 0 & \text{if } i \in I_1 \cup I_2 \cup I_3 \setminus S \\ 1 & \text{otherwise.} \end{cases}$$

Note that $f(\mathbf{x}) - f(\mathbf{x}^*) = 0.5s\kappa^2(1-4\delta) - 0.5s' \geq 0.1s\kappa^2$.

Furthermore, the gradient is equal to

$$\begin{aligned} \nabla f(\mathbf{x}) &= \mathbf{A}^\top (\mathbf{A}\mathbf{x} - \mathbf{b}) \\ &= \begin{cases} -\kappa\sqrt{1-4\delta} & \text{if } i \in I_1 \\ -\kappa\sqrt{1-2\delta} & \text{if } i \in I_2 \\ -1 & \text{if } i \in I_3 \setminus S \\ 0 & \text{if } i \in S, \end{cases} \end{aligned}$$

and since we have $\beta = \kappa$,

$$\mathbf{x} - \beta^{-1}\nabla f(\mathbf{x}) = \begin{cases} \sqrt{1-4\delta} & \text{if } i \in I_1 \\ \sqrt{1-2\delta} & \text{if } i \in I_2 \\ 1/\kappa & \text{if } i \in I_3 \setminus S \\ 1 & \text{if } i \in S, \end{cases}$$

implying that $H_{s'}(\mathbf{x} - \beta^{-1}\nabla f(\mathbf{x})) = \mathbf{x}$. □