
Sparse Convex Optimization via Adaptively Regularized Hard Thresholding

Kyriakos Axiotis¹ Maxim Sviridenko²

Abstract

The goal of *Sparse Convex Optimization* is to optimize a convex function f under a sparsity constraint $s \leq s^* \gamma$, where s^* is the target number of non-zero entries in a feasible solution (sparsity) and $\gamma \geq 1$ is an approximation factor. There has been a lot of work to analyze the sparsity guarantees of various algorithms (LASSO, Orthogonal Matching Pursuit (OMP), Iterative Hard Thresholding (IHT)) in terms of the *Restricted Condition Number* κ . The best known algorithms guarantee to find an approximate solution of value $f(x^*) + \epsilon$ with the sparsity bound of $\gamma = O\left(\kappa \min\left\{\log \frac{f(x^0) - f(x^*)}{\epsilon}, \kappa\right\}\right)$, where x^* is the target solution. We present a new *Adaptively Regularized Hard Thresholding (ARHT)* algorithm that makes significant progress on this problem by bringing the bound down to $\gamma = O(\kappa)$, which has been shown to be tight for a general class of algorithms including LASSO, OMP, and IHT. This is achieved without significant sacrifice in the runtime efficiency compared to the fastest known algorithms. We also provide a new analysis of OMP with Replacement (OMPR) for general f , under the condition $s > s^* \frac{\kappa^2}{4}$, which yields Compressed Sensing bounds under the Restricted Isometry Property (RIP). When compared to other Compressed Sensing approaches, it has the advantage of providing a strong tradeoff between the RIP condition and the solution sparsity, while working for any general function f that meets the RIP condition.

1. Introduction

Sparse Convex Optimization is the problem of optimizing a convex objective, while constraining the sparsity of the so-

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

lution (its number of non-zero entries). Variants and special cases of this problem have been studied for many years, and there have been countless applications in Machine Learning, Signal Processing, and Statistics. In Machine Learning it is used to regularize models by enforcing parameter sparsity, since a sparse set of parameters often leads to better model generalization. Furthermore, in a lot of large scale applications the number of parameters of a trained model is a significant factor in computational efficiency, thus improved sparsity can lead to improved time and memory performance. In applied statistics, a single extra feature translates to a real cost from increasing the number of samples. In Compressed Sensing, finding a sparse solution to a Linear Regression problem can be used to significantly reduce the sample size for the recovery of a target signal. In the context of these applications, decreasing sparsity by even a small amount while not increasing the accuracy can have a significant impact.

Sparse Optimization Given a function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ and any s^* -sparse (*unknown*) target solution x^* , the Sparse Optimization problem is to find an s -sparse solution x , i.e. a solution with at most s non-zero entries, such that $f(x) \leq f(x^*) + \epsilon$ and $s \leq s^* \gamma$, where $\epsilon > 0$ is a desired accuracy and $\gamma \geq 1$ is an approximation factor for the target sparsity. Even if f is a convex function, the sparsity constraint makes this problem non-convex, and it has been shown that it is an intractable problem, even when $\gamma = O\left(2^{\log^{1-\delta} n}\right)$ and f is the Linear Regression objective (Natarajan, 1995; Foster et al., 2015). However, this worst-case behavior is not observed in practice, and so a large body of work has been devoted to the analysis of algorithms under the assumption that the *restricted condition number* $\kappa_{s+s^*} = \frac{\rho_{s+s^*}^+}{\rho_{s+s^*}^-}$ (or just $\kappa = \frac{\rho^+}{\rho^-}$) of f is bounded (Natarajan, 1995; Shalev-Shwartz et al., 2010; Zhang, 2011; Bahmani et al., 2013; Liu et al., 2014; Jain et al., 2014; Yuan et al., 2016; Shen & Li, 2017a;b; Jain et al., 2014; Somani et al., 2018). Note: Here, $\rho_{s+s^*}^+$ is the maximum smoothness constant of any restriction of f on an $(s + s^*)$ -sparse subset of coordinates and $\rho_{s+s^*}^-$ is the minimum strong convexity constant of any restriction of f on an $(s + s^*)$ -sparse subset of coordinates.

The first algorithm for this problem, often called *Orthogonal Matching Pursuit (OMP)* or *Greedy*, was analyzed by

(Natarajan, 1995) for Linear Regression, and subsequently for general f by (Shalev-Shwartz et al., 2010), obtaining the guarantee that the sparsity of the returned solution is $O\left(s^* \kappa \log \frac{f(x^0) - f(x^*)}{\epsilon}\right)$ ¹. In applications where having low sparsity is crucial, the dependence of sparsity on the required accuracy ϵ is undesirable. The question of whether this dependence can be removed was answered positively (Shalev-Shwartz et al., 2010; Jain et al., 2014) giving a sparsity guarantee of $O(s^* \kappa^2)$. As remarked in (Shalev-Shwartz et al., 2010), this bound sacrifices the linear dependence on κ , while removing the dependence on ϵ and $f(x^0) - f(x^*)$.

Since then, there has been some work on improving these results by introducing non-trivial assumptions, such as the target solution x^* being close to globally optimal. More specifically, (Zhang, 2011) defines the *Restricted Gradient Optimal Constant (RGOC) at level s* , ζ_s (or just ζ) as the ℓ_2 norm of the top- s elements in $\nabla f(x^*)$ and analyzes an algorithm that gives sparsity $s = O(s^* \kappa \log(s^* \kappa))$, and such that $f(x) \leq f(x^*) + O(\zeta^2 / \rho^-)$. (Somani et al., 2018) strengthens this bound to $f(x) \leq f(x^*) + O(\zeta^2 / \rho^+)$ with sparsity $s = O(s^* \kappa \log \kappa)$. However, this means that $f(x)$ might be much larger than $f(x^*) + \epsilon$ in general. To the best of our knowledge, no improvement has been made over the $O\left(s^* \min\left\{\kappa \frac{f(x^0) - f(x^*)}{\epsilon}, \kappa^2\right\}\right)$ bound in the general case.

Another line of work studies a maximization version of the sparse convex optimization problem as well as its generalizations for matroid constraints (Altschuler et al., 2016; Elenberg et al., 2017; Chen et al., 2018).

Sparse Solution and Support Recovery Often, as is the case in Compressed Sensing, one needs a guarantee on the closeness of the solution x to the target solution x^* in absolute terms, rather than in terms of the value of f . The goal is usually either to recover (a superset of) the target support, or to ensure that the returned solution is close to the target solution in ℓ_2 norm. The results for this problem either assume a constant upper bound on the *Restricted Isometry Property (RIP)* constant $\delta_r := \frac{\kappa_r - 1}{\kappa_r + 1}$ for some r (RIP-based recovery), or that x^* is close to being a global optimum (RIP-free recovery). This problem has been extensively studied and is an active research area in the vast Compressed Sensing literature. See also the survey by (Boche et al., 2015).

In the seminal papers of (Candes & Tao, 2005; Candes et al., 2006; Donoho, 2006; Candes, 2008) it was shown that for the Linear Regression problem when $\delta_{2s^*} < \sqrt{2} - 1 \approx 0.41$, the LASSO algorithm (Tibshirani, 1996) can recover a solution with $\|x - x^*\|_2^2 \leq C f(x^*)$, where C is a constant

¹Even though (Natarajan, 1995) states a less general result, this is what is implicitly proven.

depending only on δ_{2s^*} and $f(x^*) = \frac{1}{2} \|Ax^* - b\|_2^2$ is the error of the target solution². Since then, a multitude of results of similar flavor have appeared, either giving related guarantees for the LASSO algorithm while improving the RIP upper bound (Foucart & Lai, 2009; Cai et al., 2009; Foucart, 2010; Cai et al., 2010; Mo & Li, 2011; Anderson & Strömberg, 2014) which culminate in a bound of $\delta_{2s^*} < 0.6248$, or showing that similar guarantees can be obtained by greedy algorithms under more restricted RIP conditions, but that are typically faster than LASSO (Needell & Vershynin, 2009; 2010; Needell & Tropp, 2009; Blumensath & Davies, 2009; Jain et al., 2011; Foucart, 2011; 2012). See also the comprehensive surveys (Foucart & Rauhut, 2017; Mousavi et al., 2019).

(Needell & Tropp, 2009) presents a greedy algorithm called *CoSaMP* and shows that for Linear Regression it achieves a bound in the form of (Candes, 2008) while having a more efficient implementation. Their method works for the more restricted RIP upper bound of $\delta_{2s^*} < 0.025$, or $\delta_{4s^*} < 0.4782$ as improved by (Foucart & Rauhut, 2017). (Blumensath & Davies, 2009) proves that another greedy algorithm called *Iterative Hard Thresholding (IHT)* achieves a similar bound to that of CoSaMP for Linear Regression, with the condition $\delta_{3s^*} < 0.067$, which is improved to $\delta_{2s^*} < \frac{1}{3}$ by (Jain et al., 2011) and to $\delta_{3s^*} < 0.5774$ by (Foucart, 2011).

The RIP-free line of research has shown that strong results can be achieved without a RIP upper bound, given that the target solution is sufficiently close to being a global optimum. These results typically require that s is significantly larger than s^* . In particular, (Zhang, 2011) shows that if ζ is the RGOC of f it can be guaranteed that $\|x - x^*\|_2 \leq 2\sqrt{6} \frac{\zeta}{\rho^-}$ (or $(1 + \sqrt{6}) \frac{\zeta}{\rho^-}$ with a slightly tighter analysis). (Somani et al., 2018) strengthens this bound to $\left(1 + \sqrt{1 + \frac{5}{\kappa}}\right) \frac{\zeta}{\rho^-}$. Furthermore, it has been shown that as long as a ‘‘Signal-to-Noise’’ condition holds, one can actually recover a superset of the target support. Typically the condition is a lower bound on $|x_{\min}^*|$, the minimum magnitude non-zero entry of the target solution. Different lower bounds that have been devised include $\Omega\left(\frac{\sqrt{s+s^*} \|\nabla f(x^*)\|_\infty}{\rho_{s+s^*}^-}\right)$ (Jain et al., 2014), which was later improved to $\Omega\left(\sqrt{\frac{f(x^*) - f(\bar{x}^*)}{\rho_{2s}^-}}\right)$, where \bar{x}^* is an optimal s -sparse solution (Yuan et al., 2016). Finally, (Somani et al., 2018) improves the sparsity bound to $O(s^* \kappa \log(s^* \kappa))$ in the statistical setting and (Shen & Li, 2017b) shows that the sparsity can be brought down to $s = s^* + O(\kappa^2)$ if a stronger lower bound of $\Omega\left(\sqrt{\kappa} \frac{\zeta}{\rho}\right)$ is assumed.

² $f(x^*)$ is also commonly denoted as $\frac{1}{2} \|\eta\|_2^2$, where $Ax^* = b + \eta$, i.e. η is the measurement noise.

1.1. Our work

In this work we present a new algorithm called *Adaptively Regularized Hard Thresholding (ARHT)*, that closes the longstanding gap between the $O\left(s^* \kappa \frac{f(x^0) - f(x^*)}{\epsilon}\right)$ and $O(s^* \kappa^2)$ bounds by getting a sparsity of $O(s^* \kappa)$ and thus achieving the best of both worlds. As (Foster et al., 2015) shows that for a general class of algorithms (including greedy algorithms like OMP, IHT as well as LASSO) the linear dependence on κ is necessary even for the special case of Sparse Regression, our result is tight for this class of algorithms. In the supplementary material we briefly describe this example and also state a conjecture that it can be turned into an inapproximability result. Furthermore, there we also show that the $O(s^* \kappa^2)$ sparsity bound is tight for OMPR, thus highlighting the importance of regularization in our method. Our algorithm is efficient, as it requires roughly $O\left(s \log^3 \frac{f(x^0) - f(x^*)}{\epsilon}\right)$ iterations, each of which includes one function minimization in a restricted support of size s and is simple to describe and implement. Furthermore, it directly implies non-trivial results in the area of Compressed Sensing.

We also provide a new analysis of OMPR (Jain et al., 2011) and show that under the condition that $s > s^* \frac{\kappa^2}{4}$, or equivalently under the RIP condition $\delta_{s+s^*} < \frac{2\sqrt{\frac{s}{s^*}-1}}{2\sqrt{\frac{s}{s^*}+1}}$, it is possible to approximately minimize the function f up to some error depending on the RIP constant and the closeness of x^* to global optimality. More specifically, we show that for any $\epsilon > 0$ OMPR returns a solution x such that

$$f(x) \leq f(x^*) + \epsilon + C_1(f(x^*) - f(x^{\text{opt}}))$$

where x^{opt} is the globally optimal solution, as well as

$$\|x - x^*\|_2^2 \leq \epsilon + C_2(f(x^*) - f(x^{\text{opt}}))$$

where C_1, C_2 are constants that only depend on $\frac{s}{s^*}$ and δ_{s+s^*} . An important feature of our approach is that it provides a meaningful tradeoff between the RIP constant upper bound and the sparsity of the solution, even when the sparsity s is arbitrarily close to s^* . In other words, one can relax the RIP condition at the expense of increasing the sparsity of the returned solution. Furthermore, our analysis applies to general functions with bounded RIP constant.

Experiments with real data suggest that ARHT and a variant of OMPR which we call *Exhaustive Local Search* achieve promising performance in recovering sparse solutions.

1.2. Comparison to previous work

Sparse Optimization Our Algorithm 2 (ARHT) returns a solution with $s = O(s^* \kappa)$ without any additional assumptions, thus significantly improving over the bound

Table 1. Compressed Sensing tradeoffs implied by Theorem 3.7: Sparsity vs RIP condition

s	RIP CONDITION
s^*	$\delta_{2s^*} < 0.33$
$2s^*$	$\delta_{3s^*} < 0.47$
$3s^*$	$\delta_{4s^*} < 0.55$
$30s^*$	$\delta_{31s^*} < 0.83$

$O\left(s^* \min\left\{\kappa \frac{f(x^0) - f(x^*)}{\epsilon}, \kappa^2\right\}\right)$ that was known in previous work. This proves that neither any dependence on the required solution accuracy ϵ , nor the quadratic dependence on the condition number κ is necessary. Furthermore, no assumption on the function or the target solution is required to achieve this bound. Importantly, previous results imply that our bound is tight up to constants for a general class of algorithms, including Greedy-type algorithms and LASSO (Foster et al., 2015).

Sparse Solution Recovery In Corollary 3.5, we show that the improved guarantees of Theorem 3.3 immediately imply that ARHT gives a bound of $\|x - x^*\|_2 \leq (2 + \theta) \frac{\zeta}{\rho}$ for any $\theta > 0$, where ζ is the Restricted Gradient Optimal Constant. This improves the constant factor in front of the corresponding results of (Zhang, 2011; Somani et al., 2018).

As we saw, our Theorem 3.7 directly implies that OMPR gives an upper bound on $\|x - x^*\|_2^2$ of the same form as the RIP-based bounds in previous work, under the condition $\delta_{s+s^*} < \frac{2\sqrt{\frac{s}{s^*}-1}}{2\sqrt{\frac{s}{s^*}+1}}$. While previous results either concentrate on the case $s = s^*$, or $s \gg s^*$, our analysis provides a way to trade off increased sparsity for a more relaxed RIP bound, allowing for a whole family of RIP conditions where s is arbitrarily close to s^* . Specifically, if we set $s = s^*$ our work implies recovery for $\delta_{2s^*} < \frac{1}{3} \approx 0.33$, which matches the best known bound for any greedy algorithm (Jain et al., 2011), although it is a stricter condition than the $\delta_{2s^*} < 0.62$ required by LASSO (Foucart & Rauhut, 2017). Table 1 contains a few such RIP bounds implied by our analysis and shows that it readily surpasses the bounds for Subspace Pursuit $\delta_{3s^*} < 0.35$, CoSaMP $\delta_{4s^*} < 0.48$, and OMP $\delta_{31s^*} < 0.33$ (Jain et al., 2011; Zhang, 2011). Another important feature compared to previous work is that all our guarantees are not restricted to Linear Regression and are true for any function f , as long as it satisfies the required RIP condition, which makes the result more general.

Sparse Support Recovery Corollary 3.6 shows that as a direct consequence of our work, the condition $|x_{\min}^*| > \frac{\zeta}{\rho}$ suffices for our algorithm to recover a superset of the support with size $s = O(s^* \kappa)$. Compared to (Jain et al., 2014), we improve both the size of the superset, as well as

the condition, since $\sqrt{s} \frac{\|\nabla f(x^*)\|_\infty}{\rho^-} \geq \sqrt{\frac{s}{s^*}} \frac{\zeta}{\rho^-} = \Omega\left(\frac{\zeta}{\rho^-}\right)$. Compared to (Shen & Li, 2017b), the bounds on the superset size are incomparable in general, but our $|x_{\min}^*|$ condition is more relaxed, since $\sqrt{\kappa} \frac{\zeta}{\rho^-} = \Omega\left(\frac{\zeta}{\rho^-}\right)$. Finally, compared to (Yuan et al., 2016) we have a stricter lower bound on $|x_{\min}^*|$, but with a better bound on the superset size ($O(s^* \kappa)$ instead of $O(s^* \kappa^2)$). Although not explicitly stated, (Zhang, 2011; Somani et al., 2018) also give a similar lower bound of $\sqrt{1 + \frac{10}{\kappa}} \frac{\zeta}{\rho^-}$ which we improve by a constant factor.

Runtime discussion ARHT has the advantage of being very simple to implement in practice. The runtime of Algorithm 2 (ARHT) is comparable to that of the most efficient greedy algorithms (e.g. OMP/OMPR), as it requires a single function minimization per iteration.

Naming Conventions The algorithm that we call *Orthogonal Matching Pursuit (OMP)*, is also known as ‘‘Greedy’’ (Natarajan, 1995), ‘‘Fully Corrective Forward Greedy Selection’’ or just ‘‘Forward Selection’’. What we call *Orthogonal Matching Pursuit with Replacement (OMPR)* (Jain et al., 2011) is also known by various other names. It is referenced in (Shalev-Shwartz et al., 2010) as a simpler variant of their ‘‘Fully Corrective Forward Greedy Selection with Replacement’’ algorithm, or just Forward Selection with Replacement, or ‘‘Partial Hard Thresholding with parameter $r = 1$ (PHT(r) where $r = 1$)’’ (Jain et al., 2017) which is a generalization of Iterative Hard Thresholding. Finally, what we call *Exhaustive Local Search* is essentially a variant of ‘‘Orthogonal Least Squares’’ that includes replacement steps, and also appears in (Shalev-Shwartz et al., 2010) as ‘‘Fully Corrective Forward Greedy Selection with Replacement’’, or just ‘‘Forward Stepwise Selection with Replacement’’. See also (Blumensath & Davies, 2007) regarding naming conventions.

Remark 1.1. Most of the results in the literature either only apply to, or are only presented for the Linear Regression problem. Since all of our results apply to general function minimization, we present them as such.

2. Preliminaries

Remark 2.1. An addendum to this section can be found in the *Supplementary Material*.

We denote $[i] := \{1, 2, \dots, i\}$. For any $x \in \mathbb{R}^n$ and $R \subseteq [n]$, we define $x_R \in \mathbb{R}^n$ as $(x_R)_i = \begin{cases} x_i & i \in R \\ 0 & \text{otherwise} \end{cases}$. Additionally, for any differentiable function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ with gradient $\nabla f(x)$, we will denote by $\nabla_R f(x)$ the restriction of $\nabla f(x)$ to R , i.e. $(\nabla f(x))_R$.

In Lemma 2.2 we state a standard martingale concentration

inequality that we will use. See also (Chung & Lu, 2006) for more on martingales.

Lemma 2.2 (Martingale concentration inequality (Special case of Theorem 6.5 in (Chung & Lu, 2006))). Let $Y_0 = 0, Y_1, \dots, Y_n$ be a martingale with respect to the sequence i_1, \dots, i_n such that

$$\text{Var}(Y_k \mid i_1, \dots, i_{k-1}) \leq \sigma^2$$

and

$$Y_{k-1} - Y_k \leq M$$

for all $k \in [n]$, then for any $\lambda > 0$,

$$\Pr[Y_n \leq -\lambda] \leq e^{-\lambda^2 / (2(n\sigma^2 + M\lambda/3))}$$

Definition 2.3. For any $x \in \mathbb{R}^n$, we denote the *support* of x by $\text{supp}(x) = \{i : x_i \neq 0\}$

Definition 2.4 (Restricted Condition Number). Given a differentiable function f , the *Restricted ℓ_2 -Smoothness (RSS)* constant, or just Restricted Smoothness constant, of f at sparsity level s is the minimum $\rho_s^+ \in \mathbb{R}$ such that

$$f(y) \leq f(x) + \nabla f(x)^\top (y - x) + \frac{\rho_s^+}{2} \|y - x\|_2^2$$

for all $x, y \in \mathbb{R}^n$ with $|\text{supp}(y - x)| \leq s$. Similarly, the *Restricted ℓ_2 -Strong Convexity (RSC)* constant, or just Restricted Strong Convexity constant, of f at sparsity level s is the maximum $\rho_s^- \in \mathbb{R}_+$ such that

$$f(y) \geq f(x) + \nabla f(x)^\top (y - x) + \frac{\rho_s^-}{2} \|y - x\|_2^2$$

for any $x, y \in \mathbb{R}^n$ with $|\text{supp}(y - x)| \leq s$. Given that $\rho_s^+, \rho_s^- > 0$, the *Restricted Condition Number* of f at sparsity level s is defined as $\kappa_s = \rho_s^+ / \rho_s^-$. We will also make use of $\tilde{\kappa}_s = \rho_2^+ / \rho_s^-$ which is at most κ_s as long as $s \geq 2$.

Definition 2.5 (Restricted Isometry Property (RIP)). We say that a differentiable function f has the *Restricted Isometry Property* at sparsity level s if $\rho_s^+, \rho_s^- > 0$, and the *RIP constant* of f at sparsity level s is then defined as $\delta_s = \frac{\kappa_s - 1}{\kappa_s + 1}$.

Definition 2.6 (Restricted Gradient Optimal Constant (RGOc)). Given a differentiable function f and a ‘‘target’’ solution x^* , the *Restricted Gradient Optimal Constant* (Zhang, 2011) at sparsity level s is the minimum $\zeta_s \in \mathbb{R}_+$ such that

$$|\langle \nabla f(x^*), y \rangle| \leq \zeta_s \|y\|_2$$

for all s -sparse y . Setting $y = \nabla_S f(x^*)$ for some set S with $|S| \leq s$, this implies that $\zeta_s \geq \|\nabla_S f(x^*)\|$. An alternative definition is that ζ_s is the ℓ_2 norm of the s elements of $\nabla f(x^*)$ with highest absolute value.

³We note that this is a more general definition than the one usually given for quadratic functions (i.e. Linear Regression).

Definition 2.7 (S -restricted minimizer). Given $f : \mathbb{R}^n \rightarrow \mathbb{R}$, $x^* \in \mathbb{R}^n$, and $S \subseteq [n]$, we will call x^* an S -restricted minimizer of f if $\text{supp}(x^*) \subseteq S$ and for all x such that $\text{supp}(x) \subseteq S$ we have $f(x^*) \leq f(x)$.

3. Theoretical Results

In this section we will call the current solution x and the target solution x^* , with respective supports S and S^* and sizes $s = |S|$ and $s^* = |S^*|$. The results in this section are stated in terms of $\tilde{\kappa}$, but since $\tilde{\kappa} = \frac{\rho_2^+}{\rho_{s+s^*}^-} \leq \frac{\rho_{s+s^*}^+}{\rho_{s+s^*}^-} = \kappa$, the same statements where $\tilde{\kappa}$ is replaced by κ automatically follow.

3.1. Adaptively Regularized Hard Thresholding (ARHT)

Our algorithm is essentially a Hard Thresholding algorithm (and more specifically OMPR also known as PHT(1)) with the crucial novelty that it is applied on an adaptively regularized objective function.

Hard thresholding algorithms maintain a solution x supported on $S \subseteq [n]$, which they iteratively update by inserting new elements into the support set S and removing the same number of elements from it, in order to preserve the sparsity of x . More specifically, OMPR makes one insertion and one removal in each iteration. In order to evaluate the element i to be *inserted* into S , OMPR uses the fact that, because of smoothness, $\frac{(\nabla_i f(x))^2}{2\rho_2^+}$ is a lower bound on the decrease of $f(x)$ caused by inserting i into the support, and therefore picks i to maximize $|\nabla_i f(x)|$. Similarly, in order to evaluate the element j to be *removed* from S , OMPR uses the fact that $\frac{\rho_2^+}{2} x_j^2$ upper bounds the increase of $f(x)$ caused by setting $x_j = 0$, and therefore picks j to minimize $|x_j|$. However, the real worth of j might be as small as $\frac{\rho_2^-}{2} x_j^2$, so the upper bound can be loose by a factor of $\frac{\rho_2^+}{\rho_2^-} \geq \frac{\rho_2^+}{\rho_{s+s^*}^-} = \tilde{\kappa}$.

We eliminate this discrepancy by running the algorithm on the regularized function $g(z) := f(z) + \frac{\rho_2^+}{2} \|z\|_2^2$. As the restricted condition number of g is now $O(1)$, the real worth of a removal candidate j matches the upper bound up to a constant factor.

However, even though g is now well conditioned, the analysis can only guarantee the quality of the solution in terms of the *original* objective f if the regularization is *not* applied on elements S^* , i.e. $\frac{\rho_2^+}{2} \|x_{R \setminus S^*}\|_2^2$ for some sufficiently large $R \subseteq [n]$; if this is the case, a solution with sparsity $O(s^* \tilde{\kappa})$ can be recovered. Unfortunately, there is no way of knowing a priori which elements not to regularize, as this is equivalent to finding the target solution. As a result,

the algorithm can get trapped in local minima, which are defined as states in which one iteration of the algorithm does not decrease $g(x)$, even though x is a suboptimal solution in terms of f (i.e. $f(x) > f(x^*)$).

The main contribution of this work is to characterize such local minima and devise a procedure that is able to successfully escape them, thus allowing x to converge to a desired solution for the original objective.

When x has significant ℓ_2^2 mass in the target support, the regularization term $\frac{\rho_2^+}{2} \|x\|_2^2$ might penalize the target solution too much, leading to a Type 2 iteration. In this case, we use random sampling to detect an element in the optimal support and unregularize it. This procedure escapes all local minima, thus leading to a bound in the total number of Type 2 iterations.

More concretely, we show that if at some iteration of the algorithm the value of $g(x)$ does not decrease sufficiently (Type 2 iteration), then roughly at least a $\frac{1}{\tilde{\kappa}}$ -fraction of the ℓ_2^2 mass of x lies in the target support S^* . We exploit this property by sampling an element i proportional to x_i^2 and removing its corresponding term from the regularizer (*unregularizing* it). We show that with constant probability this will happen at most $O(s^* \tilde{\kappa})$ times, as after that all the elements in S^* will have been unregularized.

The core algorithm is presented in Algorithm 1. The full algorithm additionally requires some standard routines like binary search and is presented in Algorithm 2.

Let $R \subseteq [n]$ be the set of currently regularized elements. The following invariant is a crucial ingredient for bringing the sparsity from $O(s^* \tilde{\kappa}^2)$ down to $O(s^* \tilde{\kappa})$, and we intend to enforce it at all times. It essentially states that there will always be enough elements in the current solution that are being regularized.

Invariant 3.1.

$$|R \cap S| \geq s^* \max\{1, 8\tilde{\kappa}\}$$

To give some intuition on this, ARHT owes its improved $\tilde{\kappa}$ dependence on the regularizer $\frac{\rho_2^+}{2} \|x\|_2^2$. However, during the algorithm, some elements are being unregularized. Our analysis requires that the current solution support always contains $\Omega(s^* \tilde{\kappa})$ regularized elements, which is what Invariant 3.1 states.

In the following, we will let opt denote a guess on the target value $f(x^*)$. Also, x^0 will denote the initial solution, which is an S^0 -restricted minimizer an arbitrary set $S^0 \subseteq [n]$ with $|S^0| = s$. In Algorithm 1, S^0 is defined explicitly as $[s]$, however in practice one might want to pick a better initial set (e.g. returned by running OMP).

Algorithm 1 ARHT core routine

```

1: function ARHT_core( $s, \text{opt}, \epsilon$ )
2:   function to be minimized  $f : \mathbb{R}^n \rightarrow \mathbb{R}$ 
3:   target sparsity  $s$ 
4:   target value  $\text{opt}$  (current guess for the optimal value)
5:   target error  $\epsilon$ 
6:   Define  $g_R(x) := f(x) + \frac{\rho_2^+}{2} \|x_R\|_2^2$  for all  $R \subseteq [n]$ .
7:    $R^0 \leftarrow [n]$ 
8:    $S^0 \leftarrow [s]$ 
9:    $x^0 \leftarrow \underset{\text{supp}(x) \subseteq S^0}{\text{argmin}} g_{R^0}(x)$ 
10:   $T = 2s \log \frac{g_{R^0}(x^0) - \min_x f(x)}{\epsilon}$  (number of iterations)
11:  for  $t = 0 \dots T - 1$  do
12:    if  $\min_{\text{supp}(x) \subseteq S^t} f(x) \leq \text{opt}$  then
13:      return  $\underset{\text{supp}(x) \subseteq S^t}{\text{argmin}} f(x)$ 
14:    end if
15:     $i \leftarrow \underset{i \in [n]}{\text{argmax}} |\nabla_i g_{R^t}(x^t)|$ 
16:     $j \leftarrow \underset{j \in S^t}{\text{argmin}} |x_j|$ 
17:     $S^{t+1} \leftarrow S^t \cup \{i\} \setminus \{j\}$ 
18:     $x^{t+1} \leftarrow \underset{\text{supp}(x) \subseteq S^{t+1}}{\text{argmin}} g_{R^t}(x)$ 
19:    if  $g_{R^t}(x^t) - g_{R^t}(x^{t+1}) < \frac{1}{s} (g_{R^t}(x^t) - \text{opt})$ 
20:      then
21:         $S^{t+1} \leftarrow S^t$ 
22:        Sample  $i \in R^t$  proportional to  $(x_i^t)^2$ 
23:         $R^{t+1} \leftarrow R^t \setminus \{i\}$ 
24:         $x^{t+1} \leftarrow \underset{\text{supp}(x) \subseteq S^{t+1}}{\text{argmin}} g_{R^{t+1}}(x)$ 
25:      end if
26:    end for
27:  return  $x^T$ 
end function

```

We begin with a lemma analyzing Algorithm 1, which is the core of our algorithm.

Lemma 3.2. If $s \geq s^* \max\{4\tilde{\kappa} + 7, 12\tilde{\kappa} + 6\}$ and $\text{opt} \geq f(x^*)$, with probability at least 0.2 ARHT_core(s, opt, ϵ) returns an s -sparse solution x such that $f(x) \leq \text{opt} + \epsilon$.

In other words, as long as $\text{opt} \geq f(x^*)$, a solution of value $\leq \text{opt} + \epsilon$ will be found. As the value opt is not known a priori, we perform binary search on it, as described in Algorithm 2. Furthermore, the probability of success in the previous lemma can be boosted by repeating multiple times.

The following theorem encapsulates the main result of this section.

Theorem 3.3. Given a function f and an (unknown) s^* -sparse solution x^* , with probability at least $1 - \frac{1}{n}$ Algorithm 2 returns an s -sparse solution x with $f(x) \leq f(x^*) + \epsilon$, as long as $s \geq s^* \max\{4\tilde{\kappa} + 7, 12\tilde{\kappa} + 6\}$. The number of

Algorithm 2 ARHT

```

1: function ARHT_robust( $s, \text{opt}, \epsilon, B$ )
2:   function to be minimized  $f : \mathbb{R}^n \rightarrow \mathbb{R}$ 
3:   lower bound on target value  $B$ 
4:    $x^{\text{ret}} \leftarrow \vec{0}$ 
5:   for  $z = 1 \dots 5 \log \left( 6n \log \frac{f(\vec{0}) - B}{\epsilon} \right)$  do
6:      $x \leftarrow \text{ARHT\_core}(s, \text{opt}, \epsilon)$ 
7:     if  $f(x) < f(x^{\text{ret}})$  then
8:        $x^{\text{ret}} \leftarrow x$ 
9:     end if
10:  end for
11:  return  $x^{\text{ret}}$ 
end function
13: function ARHT( $s, \epsilon$ )
14:   function to be minimized  $f : \mathbb{R}^n \rightarrow \mathbb{R}$ 
15:   target sparsity  $s$ 
16:   target error  $\epsilon$ 
17:    $B \leftarrow \min_x f(x)$ 
18:    $l \leftarrow B$ 
19:    $b \leftarrow \vec{0}$ 
20:    $r \leftarrow f(b)$ 
21:   while  $r - l > \epsilon$  do
22:      $m \leftarrow \frac{l+r}{2}$ 
23:      $x \leftarrow \text{ARHT\_robust}(s, m, \epsilon/3, B)$ 
24:     if  $f(x) > m + \epsilon/3$  then
25:        $l \leftarrow m$ 
26:     else
27:        $b \leftarrow x$ 
28:        $r \leftarrow f(x)$ 
29:     end if
30:   end while
31:   return  $b$ 
end function

```

iterations is $O \left(s \log^2 \frac{f(\vec{0}) - B}{\epsilon} \log \left(n \log \frac{f(\vec{0}) - B}{\epsilon} \right) \right)$ where $B = \min_x f(x)$.

The following corollary that bounds the total runtime can be immediately extracted. Note that in practice the total runtime heavily depends on the choice of f , and it can often be improved for various special cases (e.g. linear regression).

Corollary 3.4 (Theorem 3.3 runtime). If we denote by G the time needed to compute ∇f and by M the time to minimize f in a restricted subset of $[n]$ of size s , the total runtime of Algorithm 2 is $O \left((G + M) s \log^2 \frac{f(\vec{0}) - B}{\epsilon} \log \left(n \log \frac{f(\vec{0}) - B}{\epsilon} \right) \right)$. If gradient descent is used for the implementation of the inner optimization problem, then $M = O \left(G \tilde{\kappa} \log \frac{f(\vec{0}) - B}{\epsilon} \right)$ and so the total runtime can be bounded by $O \left(G s \tilde{\kappa} \log^3 \frac{f(\vec{0}) - B}{\epsilon} \log \left(n \log \frac{f(\vec{0}) - B}{\epsilon} \right) \right)$.

As the first corollary of the above theorem, we show that it directly implies solution recovery bounds similar to those of (Zhang, 2011), while also improving the recovery bound by a constant factor.

Corollary 3.5 (Solution recovery). Given a function f and an (unknown) s^* -sparse solution x^* , such that the Restricted Gradient Optimal Constant at sparsity level s is ζ , i.e.

$$|\langle \nabla f(x^*), y \rangle| \leq \zeta \|y\|_2$$

for all s -sparse y and as long as

$$s \geq s^* \max \{4\tilde{\kappa} + 7, 12\tilde{\kappa} + 6\}$$

Algorithm 2 ensures that

$$f(x) \leq f(x^*) + \epsilon$$

and

$$\|x - x^*\|_2 \leq \frac{\zeta}{\rho^-} \left(1 + \sqrt{1 + 2\epsilon \frac{\rho^-}{\zeta^2}} \right)$$

For any $\theta > 0$ and $\epsilon \leq \frac{\zeta^2}{\rho^-} \theta (1 + \frac{\theta}{2})$, this implies that

$$\|x - x^*\|_2 \leq (2 + \theta) \frac{\zeta}{\rho^-}$$

The next corollary shows that our Theorem 3.3 can be also used to obtain support recovery results under a ‘‘Signal-to-Noise’’ condition given as a lower bound to $|x_{\min}^*|$.

Corollary 3.6 (Support recovery). As long as

$$s \geq s^* \max \{4\tilde{\kappa} + 7, 12\tilde{\kappa} + 6\}$$

and $|x_{\min}^*| > \frac{\zeta}{\rho^-}$, Algorithm 2 with $\epsilon < -\frac{1}{2\rho^-} \zeta^2 + \frac{\rho^-}{2} (x_{\min}^*)^2$ returns a solution x with support S such that

$$S^* \subseteq S$$

3.2. Analysis of Orthogonal Matching Pursuit with Replacement (OMPR)

The OMPR algorithm was first described (under a different name) in (Shalev-Shwartz et al., 2010). It is an extension of OMP but after each iteration some element is removed from S^t so that the sparsity remains the same. It is described in Algorithm 3.

When x (with support S) and x^* (with support S^*) are clear from context, we will also define a solution

$$\tilde{x} = \operatorname{argmin}_{\operatorname{supp}(z) \subseteq S \cup S^*} f(z)$$

Algorithm 3 Orthogonal Matching Pursuit with Replacement

```

1: function OMPR( $s$ )
2:   function to be minimized  $f : \mathbb{R}^n \rightarrow \mathbb{R}$ 
3:   output sparsity  $s$ 
4:    $S^0 \leftarrow [s]$ 
5:    $x^0 \leftarrow \operatorname{argmin}\{f(x) \mid \operatorname{supp}(x) \subseteq S^0\}$ 
6:    $t \leftarrow 0$ 
7:   while true do
8:      $i \leftarrow \operatorname{argmax}\{|\nabla_i f(x^t)| \mid i \in [n] \setminus S^t\}$ 
9:      $j \leftarrow \operatorname{argmin}\{|x_j^t| \mid j \in S^t\}$ 
10:     $S^{t+1} \leftarrow S^t \cup \{i\} \setminus \{j\}$ 
11:     $x^{t+1} \leftarrow \operatorname{argmin}\{f(x) \mid \operatorname{supp}(x) \subseteq S^{t+1}\}$ 
12:    if  $f(x^{t+1}) \geq f(x^t)$  then
13:      break
14:    end if
15:     $t \leftarrow t + 1$ 
16:  end while
17:   $T \leftarrow t$ 
18:  return  $x^T$ 
19: end function

```

Furthermore, we will denote by \bar{x}^* the optimal $(s + s^*)$ -sparse solution, i.e.

$$\bar{x}^* = \operatorname{argmin}_{|\operatorname{supp}(z)| \leq s + s^*} f(z)$$

By definition, for any x, x^* the following chain of inequalities holds

$$\min_{z \in \mathbb{R}^n} f(z) \leq f(\bar{x}^*) \leq f(\tilde{x}) \leq \min\{f(x), f(x^*)\}$$

The following theorem is the main result of this section. Its strength lies in its generality, and various useful corollaries can be directly extracted from it.

Theorem 3.7. Given a function f , an (unknown) s^* -sparse solution x^* , a desired solution sparsity level s , and error parameters $\epsilon > 0$ and $0 < \theta < 1$, Algorithm 3 returns an s -sparse solution x such that

- If $\tilde{\kappa} \sqrt{\frac{s^*}{s}} \leq 1$, then

$$f(x) \leq f(x^*) + \epsilon$$

in $O\left(\sqrt{ss^*} \log \frac{f(x^0) - f(x^*)}{\epsilon}\right)$ iterations.

- If $1 < \tilde{\kappa} \sqrt{\frac{s^*}{s}} < 2 - \theta$, then

$$f(x) \leq f(x^*) + B$$

where

$$B = \epsilon + \frac{4(1 - \theta) \left(\tilde{\kappa} \sqrt{\frac{s^*}{s}} - 1 \right)}{\left(2 - \tilde{\kappa} \sqrt{\frac{s^*}{s}} - \theta \right)^2} (f(x^*) - f(\bar{x}^*))$$

in $O\left(\frac{\sqrt{ss^*}}{\theta} \log \frac{f(x^0) - f(x^*)}{B}\right)$ iterations.

The first corollary states that in the “noiseless” case (i.e. when the target solution is globally optimal), the returned solution can reach arbitrarily close to the target solution:

Corollary 3.8 (Noiseless case). If $\tilde{\kappa}\sqrt{\frac{s^*}{s}} < 2$ and x^* is a globally optimal solution, i.e. $f(x^*) = \min_z f(z)$, Algorithm 3 returns a solution with

$$f(x) \leq f(x^*) + \epsilon$$

in $O\left(\frac{\sqrt{ss^*}}{2 - \tilde{\kappa}\sqrt{\frac{s^*}{s}}} \log \frac{f(x^0) - f(x^*)}{\epsilon}\right)$ iterations.

Proof. We apply Theorem 3.7 with $\theta = \frac{1}{2}\left(2 - \tilde{\kappa}\sqrt{\frac{s^*}{s}}\right)$. \square

The following result is in the usual form of sparse recovery results, which provide a bound on $\|x - x^*\|_2$ given a RIP constant upper bound. It provides a tradeoff between the RIP constant and the sparsity of the returned solution.

Corollary 3.9 (ℓ_2 solution recovery). Given any parameters $\epsilon > 0$ and $0 < \theta < 1$, the returned solution x of Algorithm 3 will satisfy

$$\|x - x^*\|_2^2 \leq \epsilon + C \left(f(x) - \min_z f(z)\right)$$

as long as

$$\delta_{s+s^*} < \frac{(2 - \theta)\sqrt{\frac{s}{s^*}} - 1}{(2 - \theta)\sqrt{\frac{s}{s^*}} + 1}$$

where C is a constant that depends only on θ , δ_{s+s^*} , and $\frac{s}{s^*}$.

In particular, for $s = s^*$, the above lemma implies recovery under the condition $\delta_{2s^*} < \frac{1}{3}$.

4. Experiments

Remark 4.1. More details about the datasets used for evaluation can be found in the *Supplementary material*.

In this section we evaluate the training performance of different algorithms in the tasks of Linear Regression and Logistic Regression. More specifically, for each algorithm we are interested in how the *loss* over the training set (the quality of the solution) evolves as a function of the the *sparsity* of the solution, i.e. the number of non-zeros.

The algorithms that we will consider are *LASSO*, *Orthogonal Matching Pursuit (OMP)*, *Orthogonal Matching Pursuit with Replacement (OMPR)*, *Adaptively Regularized Hard*

Thresholding (ARHT) (Algorithm 2), and *Exhaustive Local Search* (which is a version of OMPR that examines all possible insertions/removals in each iteration). We run our experiments on publicly available regression and binary classification datasets, out of which we have presented those on which the algorithms have significantly different performance between each other. In some of the other datasets that we tested, we observed that all algorithms had similar performance. The results are presented in Figure 1 and Figure 2. Another relevant class of algorithms that we considered was ℓ_p *Approximate Message Passing* algorithms (Donoho et al., 2009; Zheng et al., 2017). Brief experiments showed its performance in terms of sparsity for $p \leq 0.5$ to be promising (on par with OMPR and ARHT although these had much faster runtimes), however a detailed comparison is left for future work.

In both types of objectives (linear and logistic) we include an intercept term, which is present in all solutions (i.e. it is always counted as +1 in the sparsity of the solution). For consistency, all greedy algorithms (OMPR, ARHT, Exhaustive Local Search) are initialized with the OMP solution of the same sparsity.

The experiments make it clear that Exhaustive Local Search outperforms the other algorithms. However, ARHT also has promising performance and it might be preferred because of better computational efficiency. As a general conclusion, however, both Exhaustive Local Search and ARHT offer an advantage compared to OMP and OMPR. As a limitation, we observe that ARHT has inconsistent performance in some cases, oscillating between the Exhaustive Local Search and OMPR solutions.

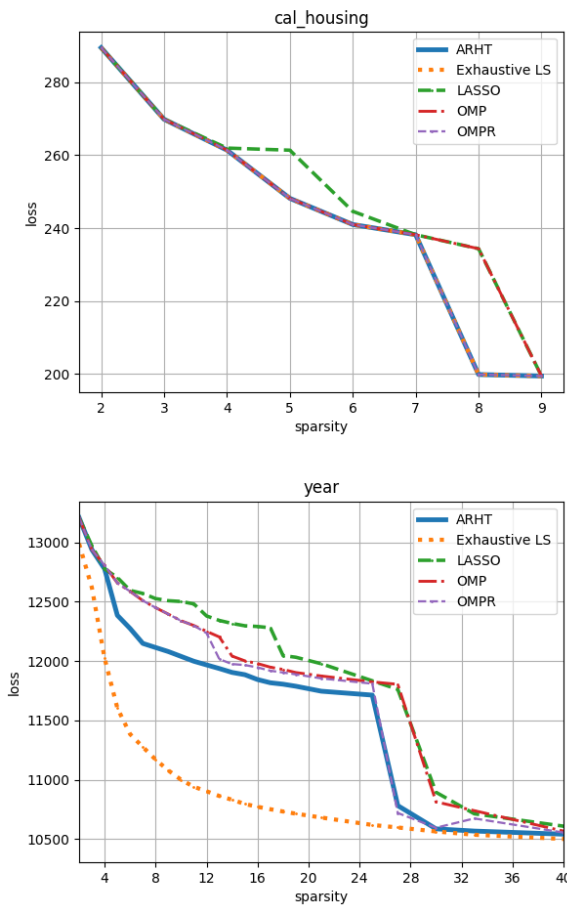


Figure 1. Comparison of different algorithms in the Regression datasets *cal_housing* and *year* using the Linear Regression loss.

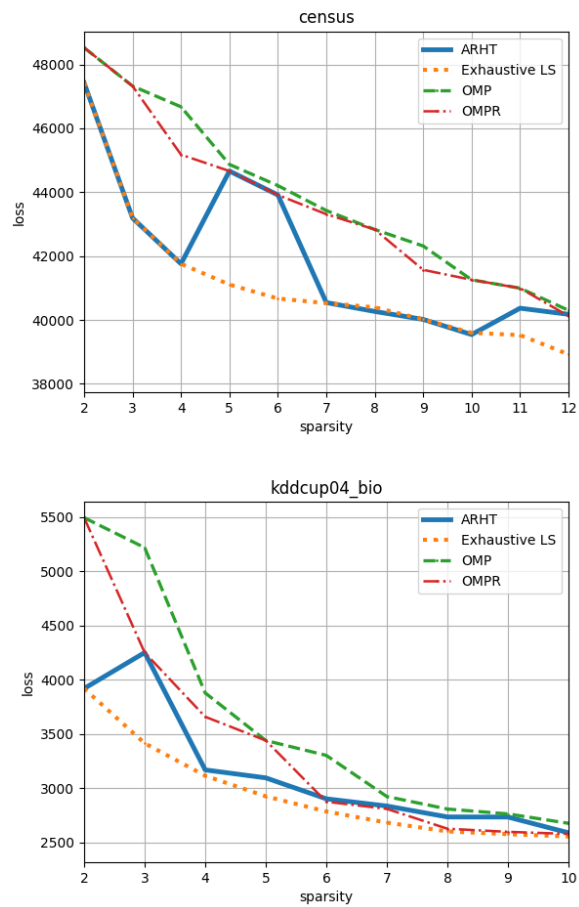


Figure 2. Comparison of different algorithms in the Binary classification datasets *census* and *kddcup04_bio* using the Logistic Regression loss.

References

- Altschuler, J., Bhaskara, A., Fu, G., Mirrokni, V., Ros-tamizadeh, A., and Zadimoghaddam, M. Greedy column subset selection: New bounds and distributed algorithms. In *International Conference on Machine Learning*, pp. 2539–2548, 2016.
- Andersson, J. and Strömberg, J.-O. On the theorem of uniform recovery of random sampling matrices. *IEEE Transactions on Information Theory*, 60(3):1700–1710, 2014.
- Bahmani, S., Raj, B., and Boufounos, P. T. Greedy sparsity-constrained optimization. *Journal of Machine Learning Research*, 14(Mar):807–841, 2013.
- Blumensath, T. and Davies, M. E. On the difference between orthogonal matching pursuit and orthogonal least squares. 2007.

- 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.
- Cai, T. T., Wang, L., and Xu, G. Shifting inequality and recovery of sparse signals. *IEEE Transactions on Signal processing*, 58(3):1300–1308, 2009.
- Cai, T. T., Wang, L., and Xu, G. New bounds for restricted isometry constants. *IEEE Transactions on Information Theory*, 56(9):4388–4394, 2010.
- Candes, E. J. The restricted isometry property and its implications for compressed sensing. *Comptes rendus mathématique*, 346(9-10):589–592, 2008.
- Candes, E. J. and Tao, T. Decoding by linear programming. *IEEE transactions on information theory*, 51(12):4203–4215, 2005.
- Candes, E. J., Romberg, J. K., and Tao, T. Stable signal recovery from incomplete and inaccurate measurements. *Communications on Pure and Applied Mathematics: A Journal Issued by the Courant Institute of Mathematical Sciences*, 59(8):1207–1223, 2006.
- Chen, L., Feldman, M., and Karbasi, A. Weakly submodular maximization beyond cardinality constraints: Does randomization help greedy? In *International Conference on Machine Learning*, pp. 804–813, 2018.
- Chung, F. and Lu, L. Concentration inequalities and martingale inequalities: a survey. *Internet Mathematics*, 3(1):79–127, 2006.
- Donoho, D. L. Compressed sensing. *IEEE Transactions on information theory*, 52(4):1289–1306, 2006.
- 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.
- Elenberg, E., Dimakis, A. G., Feldman, M., and Karbasi, A. Streaming weak submodularity: Interpreting neural networks on the fly. In *Advances in Neural Information Processing Systems*, pp. 4044–4054, 2017.
- 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. Sparse recovery algorithms: sufficient conditions in terms of restricted isometry constants. In *Approximation Theory XIII: San Antonio 2010*, pp. 65–77. Springer, 2012.
- Foucart, S. and Lai, M.-J. Sparsest solutions of underdetermined linear systems via ℓ_q -minimization for $0 < q < 1$. *Applied and Computational Harmonic Analysis*, 26(3):395–407, 2009.
- Foucart, S. and Rauhut, H. A mathematical introduction to compressive sensing. *Bull. Am. Math.*, 54:151–165, 2017.
- 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.
- Liu, J., Ye, J., and Fujimaki, R. Forward-backward greedy algorithms for general convex smooth functions over a cardinality constraint. In *International Conference on Machine Learning*, pp. 503–511, 2014.
- Mo, Q. and Li, S. New bounds on the restricted isometry constant δ_{2k} . *Applied and Computational Harmonic Analysis*, 31(3):460–468, 2011.
- Mousavi, S., Taghiabadi, M. M. R., and Ayanzadeh, R. A survey on compressive sensing: Classical results and recent advancements. *arXiv preprint arXiv:1908.01014*, 2019.
- 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.
- Needell, D. and Vershynin, R. Uniform uncertainty principle and signal recovery via regularized orthogonal matching pursuit. *Foundations of computational mathematics*, 9(3):317–334, 2009.

Needell, D. and Vershynin, R. Signal recovery from incomplete and inaccurate measurements via regularized orthogonal matching pursuit. *IEEE Journal of selected topics in signal processing*, 4(2):310–316, 2010.

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.

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.

Somani, R., Gupta, C., Jain, P., and Netrapalli, P. Support recovery for orthogonal matching pursuit: upper and lower bounds. In *Advances in Neural Information Processing Systems*, pp. 10814–10824, 2018.

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.

Zhang, T. Sparse recovery with orthogonal matching pursuit under rip. *IEEE Transactions on Information Theory*, 57(9):6215–6221, 2011.

Zheng, L., Maleki, A., Weng, H., Wang, X., and Long, T. Does ℓ_p -minimization outperform ℓ_1 -minimization? *IEEE Transactions on Information Theory*, 63(11):6896–6935, 2017.

5. Supplementary

5.1. Definitions

Definition 5.1 (ℓ_p Norms). For any $p \in \mathbb{R}_+$, we define

$$\|x\|_p = \left(\sum_i |x_i|^p \right)^{1/p}$$

as well as the special cases $\|x\|_0 = |\{i : x_i \neq 0\}|$ and $\|x\|_\infty = \max_i |x_i|$

The following lemma stems from the definitions of ρ_2^+ , ρ_1^+ and can be used to relate ρ_2^+ with ρ_1^+

Lemma 5.2. For any function f that has the RSC property at sparsity level ≥ 2 and RSS constants ρ_1^+ , ρ_2^+ at sparsity levels 1 and 2 respectively, we have $\rho_2^+ \leq 2\rho_1^+$.

Proof. For any $x, y \in \mathbb{R}^n$ such that $|\text{supp}(y - x)| \leq 2$, We will prove that

$$f(y) \leq f(x) + \langle \nabla f(x), y - x \rangle + \frac{2\rho_1^+}{2} \|y - x\|_2^2$$

Let $y = x + \alpha \vec{1}_i + \beta \vec{1}_j$ for some $i, j \in [n]$ and $\alpha, \beta \in \mathbb{R}$. We assume $i \neq j$ and since otherwise the claim already follows from RSS at sparsity level 1. We apply the RSS property with sparsity level 1 to get the inequalities

$$\begin{aligned} f(x + 2\alpha \vec{1}_i) &\leq f(x) + 2\langle \nabla f(x), \alpha \vec{1}_i \rangle + 4\frac{\rho_1^+}{2} \|\alpha \vec{1}_i\|_2^2 \end{aligned}$$

and

$$\begin{aligned} f(x + 2\beta \vec{1}_j) &\leq f(x) + 2\langle \nabla f(x), \beta \vec{1}_j \rangle + 4\frac{\rho_1^+}{2} \|\beta \vec{1}_j\|_2^2 \end{aligned}$$

Now, by using convexity (more precisely restricted convexity at sparsity level 2 that is implied by RSC) we have

$$\begin{aligned} f(y) &= f(x + \alpha \vec{1}_i + \beta \vec{1}_j) \\ &\leq \frac{1}{2} \left(f(x + 2\alpha \vec{1}_i) + f(x + 2\beta \vec{1}_j) \right) \\ &\leq f(x) + \langle \nabla f(x), \alpha \vec{1}_i + \beta \vec{1}_j \rangle + \frac{2\rho_1^+}{2} \|\alpha \vec{1}_i + \beta \vec{1}_j\|_2^2 \\ &= f(x) + \langle \nabla f(x), y - x \rangle + \frac{2\rho_1^+}{2} \|y - x\|_2^2 \end{aligned}$$

□

5.2. Algorithms

5.2.1. ℓ_1 OPTIMIZATION (LASSO)

The LASSO approach is to relax the ℓ_0 constraint into an ℓ_1 one, thus solving the following optimization problem:

$$\min_x f(x) + \lambda \|x\|_1 \quad (1)$$

for some parameter $\lambda > 0$.

5.2.2. ITERATIVE HARD THRESHOLDING (IHT):

(Blumensath & Davies, 2009) define the hard thresholding operator $H_r(x)$ as

$$[H_r(x)]_i = \begin{cases} x_i & \text{if } |x_i| \text{ is one of the } r \text{ entries of } x \\ & \text{with largest magnitude} \\ 0 & \text{otherwise} \end{cases}$$

Using this, the algorithm is described in Algorithm 4.

Algorithm 4 Iterative Hard Thresholding (IHT)

```

1: function IHT( $s, T$ )
2:   function to be minimized  $f : \mathbb{R}^n \rightarrow \mathbb{R}$ 
3:   number of iterations  $T$ 
4:   output sparsity  $s$ 
5:    $x^0 \leftarrow \vec{0}$ 
6:   for  $t = 0 \dots T - 1$  do
7:      $x^{t+1} \leftarrow H_s(x^t - \eta \nabla f(x^t))$ 
8:   end for
9:   return  $x^T$ 
10: end function
    
```

 5.2.3. ORTHOGONAL MATCHING PURSUIT
 (GREEDY/OMP/FWD STEPWISE SELECTION)

The algorithm is described in Algorithm 5.

Algorithm 5 Greedy/OMP/Fwd stepwise selection

```

1: function greedy( $s$ )
2:   function to be minimized  $f : \mathbb{R}^n \rightarrow \mathbb{R}$ 
3:   output sparsity  $s$ 
4:    $S^0 \leftarrow \emptyset$ 
5:    $x^0 \leftarrow \vec{0}$ 
6:   for  $t = 0 \dots s - 1$  do
7:      $i \leftarrow \operatorname{argmax}\{|\nabla_i f(x^t)| \mid i \in [n] \setminus S^t\}$ 
8:      $S^{t+1} \leftarrow S^t \cup \{i\}$ 
9:      $x^{t+1} \leftarrow \operatorname{argmin}\{f(x) \mid \operatorname{supp}(x) \subseteq S^{t+1}\}$ 
10:  end for
11:  return  $x^s$ 
12: end function
    
```

5.2.4. EXHAUSTIVE LOCAL SEARCH

The algorithm in this section is similar to OMPR, in that it iteratively inserts a new element in the support while removing one from it at the same time. While, as in OMPR, the element to be removed is the minimum magnitude entry, the one to be inserted is chosen to be the one which results in the maximum decrease in the value of the objective. It is described in Algorithm 6.

6. Tightness of Theorem 3.3 result for Greedy algorithms

6.1. $\Omega(s^* \kappa)$ lower bound due to (Foster et al., 2015)

In Appendix B of (Foster et al., 2015) a matrix $A \in \mathbb{R}^{m \times n}$ and a vector $b \in \mathbb{R}^m$ are constructed and let us define $f(x) = \frac{1}{2} \|Ax - b\|_2^2$. If we let $S^* = \{1, \dots, n-2\}$ and $S^* = \{n-1, n\}$, then f has the property that

$$\min_{\operatorname{supp}(x) \subseteq S^*} f(x) = \min_{\operatorname{supp}(x) \subseteq \overline{S^*}} f(x) = 0$$

Algorithm 6 Exhaustive Local Search

```

1: function to be minimized  $f : \mathbb{R}^n \rightarrow \mathbb{R}$ 
2: target sparsity  $s$ 
3: number of iterations  $T$ 
4:  $S^0 \leftarrow [s]$ 
5:  $x^0 \leftarrow \operatorname{argmin}\{f(x) \mid \operatorname{supp}(x) \subseteq S^0\}$ 
6: for  $t = 0 \dots T - 1$  do
7:    $j \leftarrow \operatorname{argmin}_{j \in S^t} x_j^2$ 
8:    $i \leftarrow \operatorname{argmin}_{i \in [n] \setminus S^t} \left\{ \min_{x : \operatorname{supp}(x) \subseteq \begin{matrix} S^t \cup \{i\} \setminus \{j\} \end{matrix}} f(x) \right\}$ 
9:    $S^{t+1} \leftarrow S^t \cup \{i\} \setminus \{j\}$ 
10:   $x^{t+1} \leftarrow \operatorname{argmin}\{f(x) \mid \operatorname{supp}(x) \subseteq S^{t+1}\}$ 
11:  if  $f(x^{t+1}) \geq f(x^t)$  then
12:    return  $x^t$ 
13:  end if
14: end for
15: return  $x^T$ 
    
```

but for any $S \subset S^*$,

$$\min_{\operatorname{supp}(x) \subseteq S} f(x) > 0$$

Furthermore, for any $S \subset \overline{S^*}$ and $x = \operatorname{argmin}_{\operatorname{supp}(x) \subseteq S} f(x)$, it is true that

$$\max_{i \in S^*} |\nabla_i f(x)| < \min_{i \in \overline{S^*} \setminus S} |\nabla_i f(x)|$$

This means that for any algorithm with an OMP-like criterion like Orthogonal Matching Pursuit, Orthogonal Matching Pursuit with Replacement, Iterative Hard Thresholding, and Partial Hard Thresholding, if the initial solution does not have an intersection with S^* , then it will never have, therefore implying that the sparsity returned by the algorithm is $|S| = n - 2 = \Omega(n)$. As for this construction $\kappa = \frac{\rho_n^+}{\rho_n} = O(n)$, there exists a constant c such that the sparsity of the returned solution cannot be less than $cs^* \kappa$, since $s^* \kappa = O(n) = O(|S|)$. Therefore none of these algorithms can improve the bound $O(s^* \kappa)$ of Theorem 3.3 by more than a constant factor. This example also applies to ARHT and Exhaustive Local Search.

It seems difficult to get past this example and achieve sparsity $s = O(s^* \kappa^{1-\delta})$ for some $\delta > 0$. We conjecture that there might be a way to turn the above example into an inapproximability result:

Conjecture 6.1. For any $\delta > 0$, there is no polynomial time algorithm that given a matrix $A \in \mathbb{R}^{m \times n}$, a vector $b \in \mathbb{R}^m$, a target sparsity $s^* \geq 1$, and a desired accuracy $\epsilon > 0$, returns an $s = O(s^* \kappa_{s+s^*}^{1-\delta})$ -sparse solution x such that

$\|Ax - b\|_2^2 \leq \min_{\|x^*\|_0 \leq s^*} \|Ax^* - b\|_2^2 + \epsilon$, if such a solution exists.

6.2. $\Omega(s^* \kappa^2)$ lower bound for OMPR

The following lemma shows that, without regularization, OMPR requires sparsity $\Omega(s^* \kappa^2)$ in general, and therefore the sparsity upper bound is tight. We assume that the algorithm is run for T iterations (even when the solution stops improving).

Lemma 6.2. There is a function $f(x) = \frac{1}{2} \|Ax - b\|_2^2$ where $A \in \mathbb{R}^{n \times n}$ and $b \in \mathbb{R}^n$ and a target solution x^* of f with sparsity s^* , as well as a set $S \subseteq [n]$ with $|S| = \Theta(s^* \kappa^2)$ such that OMPR initialized with support set S returns a solution x with $f(x) = f(x^*) + \Theta(s^* \kappa^2)$.

Proof. Without loss of generality we assume that κ is an even integer and set $n = s^* (1 + \kappa + \kappa^2)$. We then partition $[n]$ into three intervals $I_1 = [1, s^*]$, $I_2 = [s^* + 1, s^*(1 + \kappa)]$, $I_3 = [s^*(1 + \kappa) + 1, s^*(1 + \kappa + \kappa^2)]$. We define the diagonal matrix $A \in \mathbb{R}^{n \times n}$ such that

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

and vector $b \in \mathbb{R}^n$ such that

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

where $\delta > 0$ is a sufficiently small scalar used to avoid ties in the steps of the algorithm. The target solution is defined as

$$x_i^* = \begin{cases} \kappa(1 - 4\delta) & \text{if } i \in I_1 \\ 0 & \text{if } i \in I_2 \cup I_3 \end{cases}$$

and its value is $f(x^*) = s^* \kappa^2 (1 - \delta)$. Now consider any initial support set $S^0 \subset I_3$ such that $|S^0| = s^* \kappa^2 / 2$. The initial solution will then be

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

and its value $f(x^0) = s^* \kappa^2 (\frac{5}{4} - 3\delta) = f(x^*) + \Theta(s^* \kappa^2)$. The gradient at x^0 is

$$\nabla_i f(x^0) = \begin{cases} -\kappa \sqrt{1 - 4\delta} & \text{if } i \in I_1 \\ -\kappa \sqrt{1 - 2\delta} & \text{if } i \in I_2 \\ -1 & \text{if } i \in I_3 \setminus S^0 \\ 0 & \text{if } i \in S^0 \end{cases}$$

therefore the algorithm will pick $S^1 = S^0 \cup \{i^0\} \setminus \{j^0\}$ for some $i^0 \in I_2$ and some $j^0 \in S^0$, since the gradient entries in I_2 have the largest magnitude among those in $[n]$. The new solution will be

$$x_i^1 = \begin{cases} 0 & \text{if } i \in I_1 \cup I_2 \cup I_3 \setminus S^1 \\ \sqrt{1 - 2\delta} & \text{if } i = i^0 \\ 1 & \text{if } i \in S^1 \setminus \{i^0\} \end{cases}$$

with value $f(x^1) = s^* \kappa^2 (\frac{5}{4} - 3\delta) - \frac{1}{2} (\kappa(1 - 2\delta) - 1)$ and gradient

$$\nabla_i f(x^1) = \begin{cases} -\kappa \sqrt{1 - 4\delta} & \text{if } i \in I_1 \\ -\kappa \sqrt{1 - 2\delta} & \text{if } i \in I_2 \setminus S^1 \\ -1 & \text{if } i \in I_3 \setminus S^1 \\ 0 & \text{if } i \in S^1 \end{cases}$$

and therefore the algorithm will pick $S^2 = S^1 \cup \{i^1\} \setminus \{i^0\}$ for some $i^1 \in I_2$. i^0 will be the one to be removed from S^1 because x_{i^0} has the smallest magnitude out of all entries in S^1 . Continuing this process, the algorithm will always have $S^t \cap I_2 = 1$ and $S^t \cap I_3 = |S^t| - 1$, and so $f(x^t) = s^* \kappa^2 (\frac{5}{4} - 3\delta) - \frac{1}{2} (\kappa(1 - 2\delta) - 1) = f(x^*) + \Theta(s^* \kappa^2)$ for $t \geq 1$. \square

7. Proofs

In the following, we will denote the minimization objective by $f(x)$. We will also write its RSS and RSC parameters ρ_2^+ and $\rho_{s+s^*}^-$ as ρ^+ and ρ^- respectively, as well as $\tilde{\kappa} = \frac{\rho_2^+}{\rho_{s+s^*}^-}$.

The use of ρ_2^+ is not restrictive. As shown in Lemma 5.2, $\rho_2^+ \leq 2\rho_1^+$ and so in all the bounds involving $\tilde{\kappa}$, it can be replaced by $2\frac{\rho_1^+}{\rho_{s+s^*}^-}$, therefore only losing a factor of 2.

7.1. Adaptively Regularized Hard Thresholding (ARHT)

We will denote the regularized objective by $g^t(x) = f(x) + \Phi^t(x)$. Let s be such that $s \geq s^* \max\{4\tilde{\kappa} + 7, 12\tilde{\kappa} + 6\}$ and $(\rho_2^+)'$ and $(\rho_{s+s^*}^-)'$ be RSS and RSC parameters of g^t . Now, for some $R^t \subseteq [n]$, let $\Phi^t(x) = \frac{\rho_2^+}{2} \|x_{R^t}\|_2^2$ be the regularizer. We start with a lemma regarding the RSC and RSS of the regularized function.

Lemma 7.1 (RSC, RSS of regularized function). $(\rho_2^+)' \leq 2\rho_2^+$ and $(\rho_{s+s^*}^-)' \geq \rho_{s+s^*}^-$

Proof. Φ^t is a quadratic restricted on R^t

$$\begin{aligned} & \Phi(y) - \Phi(x) - \nabla\Phi(x)^T(y-x) \\ &= \frac{\rho_2^+}{2} \left(\|y_{R^t}\|_2^2 - \|x_{R^t}\|_2^2 - 2x_{R^t}^T(y_{R^t} - x_{R^t}) \right) \\ &= \frac{\rho_2^+}{2} \|y_{R^t} - x_{R^t}\|_2^2 \in \left[0, \frac{\rho_2^+}{2} \|y - x\|_2^2 \right] \end{aligned}$$

and so for any x, y with $|\text{supp}(y-x)| \leq s + s^*$ (resp. $|\text{supp}(y-x)| \leq 1$) we have

$$\begin{aligned} & g(y) - g(x) - \nabla g(x)^T(y-x) \\ &= f(y) - f(x) - \nabla f(x)^T(y-x) \\ &+ \Phi(y) - \Phi(x) - \nabla\Phi(x)^T(y-x) \\ &\geq \frac{\rho_{s+s^*}^-}{2} \|y-x\|_2^2 \text{ (resp. } \leq \rho_2^+ \|y-x\|_2^2) \end{aligned}$$

□

Now, when referring to Algorithm 1, let our current solution be x^t , where x^t is an S^t -restricted minimizer and the optimal solution be x^* with support $S^* = \text{supp}(x^*)$. We will use the shorthand $\rho^+ = \rho_2^+$ and $\rho^- = \rho_{s+s^*}^-$. Note that by definition of the algorithm, x^t is an S^t -restricted minimizer of g^t . When referring to Algorithm 1, let us call iterations *Type 1* if the condition in Line 19 is false, and *Type 2* otherwise.

Lemma 7.2. If $\text{opt} \geq f(x^*)$, the progress $g^t(x^t) - g^t(x^{t+1})$ in Line 19 of Algorithm 1 is

$$\begin{aligned} & g^t(x^t) - g^t(x^{t+1}) \\ &\geq \frac{\rho^-}{2|S^* \setminus S^t| \rho^+} \left(f(x^t) - f(x^*) + \langle \nabla_{S^t \setminus S^*} \Phi^t(x^t), x_{S^t \setminus S^*}^t \rangle \right. \\ &\quad \left. - \frac{1}{2\rho^-} \|\nabla_{S^t \cap S^*} \Phi^t(x^t)\|_2^2 \right) - \rho^+(x_j^t)^2 \end{aligned}$$

Proof. First of all, since the condition in Line 12 (“if $\min_{x \subseteq S^t} f(x) \leq \text{opt}$ ”) was not triggered, we have that $\min_{x \subseteq S^t} f(x) > \text{opt} \geq f(x^*)$ and so $S^* \setminus S^t \neq \emptyset$. By Lemma 7.1 we have that $(\rho^+)' \leq 2\rho^+$, therefore the decrease in g^t that is achieved is

$$\begin{aligned} & g^t(x^t) - g^t(x^{t+1}) \\ &\geq \max_{\eta \in \mathbb{R}} \left\{ g^t(x^t) - g^t(x^t + \eta \vec{1}_i - x_j^t \vec{1}_j) \right\} \\ &\geq \max_{\eta \in \mathbb{R}} \left\{ -\langle \nabla g^t(x^t), \eta \vec{1}_i - x_j^t \vec{1}_j \rangle - \rho^+ \eta^2 - \rho^+(x_j^t)^2 \right\} \\ &:= B \end{aligned}$$

Note that, as defined by the algorithm, x^t is an S^t -restricted minimizer of g^t and since $j \in S^t$, we have $\nabla_j g^t(x^t) = 0$. Therefore

$$\begin{aligned} B &= \max_{\eta} \left\{ -\langle \nabla g^t(x^t), \eta \vec{1}_i \rangle - \rho^+ \eta^2 - \rho^+(x_j^t)^2 \right\} \\ &= \frac{[\nabla_i g^t(x^t)]^2}{4\rho^+} - \rho^+(x_j^t)^2 \\ &\geq \max_{k \in S^* \setminus S^t} \frac{[\nabla_k g^t(x^t)]^2}{4\rho^+} - \rho^+(x_j^t)^2 \\ &\geq \frac{\|\nabla_{S^* \setminus S^t} g^t(x^t)\|_2^2}{4|S^* \setminus S^t| \rho^+} - \rho^+(x_j^t)^2 \end{aligned} \tag{2}$$

where we used the fact that i was picked to maximize $|\nabla_k g^t(x^t)|$. Now we would like to relate this to $g^t(x^t) - f(x^*)$ (and not $g^t(x^t) - g^t(x^*)$). By applying the Restricted Strong Convexity property,

$$\begin{aligned} & f(x^*) - f(x^t) \\ &\geq \langle \nabla f(x^t), x^* - x^t \rangle + \frac{\rho^-}{2} \|x^t - x^*\|_2^2 \\ &\geq \langle \nabla f(x^t), x^* - x^t \rangle \\ &\quad + \frac{\rho^-}{2} \|x_{S^* \setminus S^t}^*\|_2^2 + \frac{\rho^-}{2} \|(x^t - x^*)_{S^t \cap S^*}\|_2^2 \end{aligned}$$

Now note that $f(x^t) = g^t(x^t) - \Phi^t(x^t)$, $\nabla_{S^t} g^t(x^t) = \vec{0}$ (since x^t is an S^t -restricted minimizer of g^t), and $\nabla \Phi^t(x^t) = \nabla_{S^t} \Phi^t(x^t)$ therefore

$$\begin{aligned} & \langle \nabla f(x^t), x^* - x^t \rangle \\ &= \langle \nabla g^t(x^t), x^* - x^t \rangle - \langle \nabla \Phi^t(x^t), x^* - x^t \rangle \\ &= \langle \nabla_{S^* \setminus S^t} g^t(x^t), x_{S^* \setminus S^t}^* \rangle + \langle \nabla_{S^t \setminus S^*} \Phi^t(x^t), x_{S^t \setminus S^*}^t \rangle \\ &\quad + \langle \nabla_{S^t \cap S^*} \Phi^t(x^t), (x^t - x^*)_{S^t \cap S^*} \rangle \end{aligned}$$

Plugging this into the previous inequality, we get

$$\begin{aligned} & f(x^*) - f(x^t) \\ &\geq \langle \nabla_{S^* \setminus S^t} g^t(x^t), x_{S^* \setminus S^t}^* \rangle + \frac{\rho^-}{2} \|x_{S^* \setminus S^t}^*\|_2^2 \\ &\quad + \langle \nabla_{S^t \setminus S^*} \Phi^t(x^t), x_{S^t \setminus S^*}^t \rangle \\ &\quad + \langle \nabla_{S^t \cap S^*} \Phi^t(x^t), (x^t - x^*)_{S^t \cap S^*} \rangle + \frac{\rho^-}{2} \|(x^t - x^*)_{S^t \cap S^*}\|_2^2 \\ &\geq -\frac{1}{2\rho^-} \|\nabla_{S^* \setminus S^t} g^t(x^t)\|_2^2 + \langle \nabla_{S^t \setminus S^*} \Phi^t(x^t), x_{S^t \setminus S^*}^t \rangle \\ &\quad - \frac{1}{2\rho^-} \|\nabla_{S^t \cap S^*} \Phi^t(x^t)\|_2^2 \end{aligned}$$

where we twice used the inequality $\langle u, v \rangle + \frac{\lambda}{2} \|v\|_2^2 \geq -\frac{1}{2\lambda} \|u\|_2^2$ for any $\lambda > 0$. This inequality is derived by expanding $\frac{1}{2} \left\| \frac{1}{\sqrt{\lambda}} u + \sqrt{\lambda} v \right\|_2^2 \geq 0$. So plugging in

$\|\nabla_{S^* \setminus S^t} g^t(x^t)\|_2^2$ to (2),

$$B \geq \frac{\rho^-}{2|S^* \setminus S^t| \rho^+} \left(f(x^t) - f(x^*) + \langle \nabla_{S^t \setminus S^*} \Phi^t(x^t), x_{S^t \setminus S^*}^t \rangle - \frac{1}{2\rho^-} \|\nabla_{S^t \cap S^*} \Phi^t(x^t)\|_2^2 \right) - \rho^+(x_j^t)^2$$

□

In the following we will use the following important property of our regularizer:

Observation 7.3. The definition of $\Phi^t(x^t)$ implies that

$$\begin{aligned} \langle \nabla_{S^t \setminus S^*} \Phi^t(x^t), x_{S^t \setminus S^*}^t \rangle &= \rho^+ \langle x_{R^t \setminus S^*}^t, x_{S^t \setminus S^*}^t \rangle \\ &= \rho^+ \left\| x_{R^t \setminus S^*}^t \right\|_2^2 \end{aligned}$$

and

$$\left\| \nabla_{S^t \cap S^*} \Phi^t(x^t) \right\|_2^2 = (\rho^+)^2 \left\| x_{R^t \cap S^*}^t \right\|_2^2$$

which also means that

$$\begin{aligned} \Phi^t(x^t) &= \frac{1}{2} \langle \nabla_{S^t \setminus S^*} \Phi^t(x^t), x_{S^t \setminus S^*}^t \rangle + \frac{1}{2\rho^+} \left\| \nabla_{S^t \cap S^*} \Phi^t(x^t) \right\|_2^2 \end{aligned}$$

We proceed to show that with constant probability Algorithm 1 will only have $O(s^* \tilde{\kappa})$ Type 2 iterations.

Lemma 7.4. If $\text{opt} \geq f(x^*)$, then with probability at least 0.2 the number of Type 2 iterations is at most $(s^* - 1)(4\tilde{\kappa} + 6)$.

Proof. Initially we have $|R^0 \cap S^0| = |S^0| = s$. Since in each Type 2 iteration we have $R^{t+1} = R^t - 1$,

$$|R^t \cap S^t| \geq s - [\text{number of Type 2 iterations up to } t]$$

This implies that for the first $(s^* - 1)(4\tilde{\kappa} + 6)$ Type 2 iterations,

$$|R^t \cap S^t| \geq s - (s^* - 1)(4\tilde{\kappa} + 6) \geq s^* \max\{1, 8\tilde{\kappa}\} \quad (3)$$

since $s \geq s^* \max\{4\tilde{\kappa} + 7, 12\tilde{\kappa} + 6\}$. From this it follows that

$$\begin{aligned} |(R^t \cap S^t) \setminus S^*| &= |R^t \cap S^t| - |R^t \cap S^t \cap S^*| \\ &\geq s^* \max\{1, 8\tilde{\kappa}\} - |S^t \cap S^*| \\ &\geq |S^* \setminus S^t| 8\tilde{\kappa} \\ &= |S^* \setminus S^t| 8 \frac{\rho^+}{\rho^-} \end{aligned}$$

and so

$$\begin{aligned} (x_j^t)^2 &\leq \frac{1}{|(R^t \cap S^t) \setminus S^*|} \left\| x_{(R^t \cap S^t) \setminus S^*}^t \right\|_2^2 \\ &\leq \frac{\rho^-}{8|S^* \setminus S^t| \rho^+} \left\| x_{R^t \setminus S^*}^t \right\|_2^2 \\ &= \frac{\rho^-}{8|S^* \setminus S^t| (\rho^+)^2} \langle \nabla_{S^t \setminus S^*} \Phi^t(x^t), x_{S^t \setminus S^*}^t \rangle \end{aligned}$$

by Observation 7.3. Combining this inequality with the statement of Lemma 7.2 we have

$$\begin{aligned} g^t(x^t) - g^t(x^{t+1}) &\geq \frac{\rho^-}{2|S^* \setminus S^t| \rho^+} \left(f(x^t) - f(x^*) + \langle \nabla_{S^t \setminus S^*} \Phi^t(x^t), x_{S^t \setminus S^*}^t \rangle - \frac{1}{2\rho^-} \|\nabla_{S^t \cap S^*} \Phi^t(x^t)\|_2^2 \right) - \rho^+(x_j^t)^2 \\ &\geq \frac{\rho^-}{2|S^* \setminus S^t| \rho^+} \left(f(x^t) - f(x^*) + \frac{3}{4} \langle \nabla_{S^t \setminus S^*} \Phi^t(x^t), x_{S^t \setminus S^*}^t \rangle - \frac{1}{2\rho^-} \|\nabla_{S^t \cap S^*} \Phi^t(x^t)\|_2^2 \right) \end{aligned} \quad (4)$$

By definition of a Type 2 iteration,

$$\begin{aligned} g^t(x^t) - g^t(x^{t+1}) &< \frac{1}{s} (g^t(x^t) - \text{opt}) \\ &\leq \frac{\rho^-}{2|S^* \setminus S^t| \rho^+} (g^t(x^t) - f(x^*)) \\ &= \frac{\rho^-}{2|S^* \setminus S^t| \rho^+} (f(x^t) - f(x^*) + \Phi^t(x^t)) \end{aligned} \quad (5)$$

where we used the fact that $s \geq 2s^* \tilde{\kappa} \geq 2|S^* \setminus S^t| \tilde{\kappa}$ and $f(x^*) \leq \text{opt}$. Combining inequalities (4) and (5) we get

$$\begin{aligned} \Phi^t(x^t) &> \frac{3}{4} \langle \nabla_{S^t \setminus S^*} \Phi^t(x^t), x_{S^t \setminus S^*}^t \rangle - \frac{1}{2\rho^-} \left\| \nabla_{S^t \cap S^*} \Phi^t(x^t) \right\|_2^2 \end{aligned}$$

or equivalently, by replacing $\Phi^t(x^t)$ from Observation 7.3,

$$\begin{aligned} &\frac{1}{2} \left(\frac{1}{\rho^-} + \frac{1}{\rho^+} \right) \left\| \nabla_{S^t \cap S^*} \Phi^t(x^t) \right\|_2^2 \\ &> \frac{1}{4} \langle \nabla_{S^t \setminus S^*} \Phi^t(x^t), x_{S^t \setminus S^*}^t \rangle \end{aligned}$$

Further applying Observation 7.3, we equivalently get

$$2(1 + \tilde{\kappa}) \left\| x_{R^t \cap S^*}^t \right\|_2^2 > \left\| x_{R^t \setminus S^*}^t \right\|_2^2 \quad (6)$$

Now, note that in Lines 21-22 the algorithm picks an element $i \in R^t$ with probability proportional to $(x_i^t)^2$ and unregularizes it, i.e. sets $R^{t+1} \leftarrow R^t \setminus \{i\}$. We denote this probability distribution over $i \in R^t$ by \mathcal{D} . From what we

have established already in (6), we can lower bound the probability that i lies in the target support:

$$\begin{aligned} \Pr_{i \sim \mathcal{D}}[i \in S^*] &= \frac{\|x_{R^t \cap S^*}^t\|_2^2}{\|x_{R^t \cap S^*}^t\|_2^2 + \|x_{R^t \setminus S^*}^t\|_2^2} \\ &> \frac{\frac{1}{2(1+\tilde{\kappa})}}{1 + \frac{1}{2(1+\tilde{\kappa})}} \\ &= \frac{1}{2\tilde{\kappa} + 3} \\ &:= p \end{aligned} \quad (7)$$

Note that this event can happen at most once for each $i \in S^*$ during the whole execution of the algorithm, since each element can only be removed once from the set of regularized elements.

We will prove that with constant probability the number of Type 2 steps will be at most $(s^* - 1)(4\tilde{\kappa} + 6) := b$. For $1 \leq k \leq b$, we define the following random variables:

- $i_k \in [n]$ is the index picked in the k -th Type 2 iteration, or \perp if there are less than k Type 2 iterations.
- q_k is the probability of picking an index in the optimal support in the k -th Type 2 iteration (i.e. $i_k \in S^*$):

$$q_k = \begin{cases} \|x_{R^{t_k} \cap S^*}^{t_k}\|_2^2 / \|x_{R^{t_k}}^{t_k}\|_2^2 & \text{if } i_k \neq \perp \\ 0 & \text{otherwise} \end{cases}$$

where $t_k \in [T]$ is the index of the k -th Type 2 iteration within all iterations of the algorithm.

- X_k is 1 if the index picked in the k -th Type 2 step was in the optimal support:

$$X_k = \begin{cases} 1 & \text{with probability } q_k \\ 0 & \text{otherwise} \end{cases}$$

Our goal is to upper bound $\Pr\left[\sum_{k=1}^b X_k \leq s^* - 1\right]$. This automatically implies the same upper bound on the probability that there will be more than b Type 2 iterations.

We define another sequence of random variables Y_0, \dots, Y_b , where $Y_0 = 0$, and

$$Y_k = \begin{cases} Y_{k-1} + \frac{p}{q_k} - p & \text{if } X_k = 1 \\ Y_{k-1} - p & \text{if } X_k = 0 \end{cases}$$

for $k \in [b]$. Since if $q_k > 0$ we have $\frac{p}{q_k} \leq 1$, it is immediate that

$$Y_k - Y_{k-1} \leq X_k - p$$

and so $Y_b \leq \sum_{k=1}^b X_k - bp$. Furthermore,

$$\begin{aligned} \mathbb{E}[Y_k | i_1, \dots, i_{k-1}] \\ &= Y_{k-1} + q_k \left(\frac{p}{q_k} - p \right) - (1 - q_k)p \\ &= Y_{k-1} \end{aligned}$$

meaning that Y_0, \dots, Y_b is a martingale with respect to i_1, \dots, i_b . We will apply the inequality from Lemma 2.2. We compute a bound on the differences

$$\begin{aligned} Y_{k-1} - Y_k \\ &= \begin{cases} p - \frac{p}{q_k} & \text{if } X_k = 1 \\ p & \text{if } X_k = 0 \end{cases} \\ &\leq p \end{aligned} \quad (8)$$

and the variance

$$\begin{aligned} \text{Var}(Y_k | i_1, \dots, i_{k-1}) \\ &= \mathbb{E}\left[(Y_k - \mathbb{E}[Y_k | i_1, \dots, i_{k-1}])^2 | i_1, \dots, i_{k-1}\right] \\ &= \mathbb{E}\left[(Y_k - Y_{k-1})^2 | i_1, \dots, i_{k-1}\right] \\ &= q_k \cdot \left(p - \frac{p}{q_k}\right)^2 + (1 - q_k) \cdot p^2 \\ &= q_k \cdot p^2 \left(1 - \frac{2}{q_k} + \frac{1}{q_k^2}\right) + (1 - q_k) \cdot p^2 \\ &= p^2 \left(\frac{1}{q_k} - 1\right) \\ &\leq p \end{aligned}$$

where we used (8) along with the fact that $q_k \geq p$. Using the concentration inequality from Lemma 2.2 we obtain

$$\begin{aligned} \Pr\left[\sum_{k=1}^b X_k \leq s^* - 1\right] \\ &\leq \Pr[Y_b \leq s^* - 1 - b \cdot p] \\ &\leq e^{-(bp - s^* + 1)^2 / (2(b \cdot p + p \cdot (bp - s^* + 1)/3))} \\ &= e^{-(s^* - 1) / (2(2 + p/3))} \\ &\leq e^{-1 / (2(2 + 1/9))} \\ &< 0.8 \end{aligned}$$

where we used the fact that $bp = 2(s^* - 1)$, $s^* \geq 2$ (otherwise the problem is trivial), and $p = \frac{1}{2\tilde{\kappa} + 3} \leq \frac{1}{3}$. Therefore we conclude that the probability that we have not unregularized the whole set S^* after b steps is at most 0.8. Since we can only have a Type 2 step if there is a regularized element in S^* (this is immediate from (7)), this implies that with probability at least 0.2 the number of Type 2 steps is at most $b = (s^* - 1)(4\tilde{\kappa} + 6)$. \square

Proof of Lemma 3.2.

By Lemma 7.4, with probability at least 0.2 there will be at most $(s^* - 1)(4\tilde{\kappa} + 6)$ Type 2 iterations. This means that the number of Type 1 iterations is at least $T - (s^* - 1)(4\tilde{\kappa} + 6) \geq s \log \frac{g^0(x^0) - f(x^*)}{\epsilon}$. Lemma 7.7 then implies that $f(x^T) \leq f(x^*) + \epsilon$. \square

We turn the result of Lemma 3.2 into a high probability result by repeating multiple times:

Lemma 7.5. As long as $s \geq s^* \max\{4\tilde{\kappa} + 7, 12\tilde{\kappa} + 6\}$ and $\text{opt} \geq f(x^*)$, $\text{ARHT_robust}(s, \text{opt}, \epsilon, B)$ returns an s -sparse solution x such that $f(x) \leq \text{opt} + \epsilon$ with probability at least $1 - \frac{1}{6n \log \frac{f(\bar{0}) - B}{\epsilon}}$.

Proof. From Lemma 3.2, the probability that a given call to ARHT_core fails is at most 0.8. Since this random experiment is executed $5 \log \left(6n \log \frac{f(\bar{0}) - B}{\epsilon}\right)$ times independently, the probability that it never succeeds is at most $(0.8)^{5 \log \left(6n \log \frac{f(\bar{0}) - B}{\epsilon}\right)} < \frac{1}{6n \log \frac{f(\bar{0}) - B}{\epsilon}}$, therefore the statement follows. \square

Lemma 7.6. As long as $s \geq s^* \max\{4\tilde{\kappa} + 7, 12\tilde{\kappa} + 6\}$, $\text{ARHT}(s, \epsilon)$ (in Algorithm 2) returns an s -sparse solution x such that $f(x) \leq f(x^*) + \epsilon$. The algorithm succeeds with probability at least $1 - \frac{1}{n}$, and the number of calls to ARHT_robust is $\leq 6 \log \frac{f(\bar{0}) - B}{\epsilon}$.

Proof. First we will bound the number of calls to ARHT_robust . Let L_k be the equal to $r - l$ before the k -th iteration in Line 21 of Algorithm 2. Then either $L_{k+1} = L_k/2$ (Line 25) or $L_{k+1} \leq L_k/2 + \epsilon/3 < 5L_k/6$ (Line 28). Therefore in any case we have $L_{k+1} < 5L_k/6$ which implies that after $T = 6 \log \frac{f(\bar{0}) - B}{\epsilon}$ iterations we will have $r - l \leq \epsilon$.

Now let us compute the probability that all the calls to ARHT_robust are successful. The number of such calls is at most $6 \log \frac{f(\bar{0}) - B}{\epsilon}$ and we know each one of them independently fails with probability less than $\frac{1}{6n \log \frac{f(\bar{0}) - f(B)}{\epsilon}}$, so by a union bound the probability that at least one call fails is less than $\frac{1}{n}$.

To prove correctness, note that by Lemma 7.5, for each $r \geq f(x^*)$ we have $f(\text{ARHT_robust}(s, r, \epsilon/3, B)) \leq r + \epsilon/3$. After Line 20 of Algorithm 2, we will have $l = B \leq f(x^*)$. In the while construct, it is always true that $f(x^*) \geq l$. This is initially true, as we saw. For each m chosen in Line 22 and x in Line 23, note that if $f(x) > m + \epsilon/3$, then by Lemma 7.5 $f(x^*) > m$ and so the invariant that $f(x^*) \geq l$ stays true. On the other hand, it is always true that $f(b) \leq r$. Initially this is so because $f(\bar{0}) = r$, and when we decrease r to some $f(x)$ we also update $b = x$.

This implies that in the end of the algorithm the returned solution will have the required property, since we will have $f(b) \leq r \leq l + \epsilon \leq f(x^*) + \epsilon$. \square

Proof of Theorem 3.3. Lemma 7.6 already establishes the correctness of the algorithm whp. For the runtime, note that ARHT_core takes $O\left(s \log \frac{f(\bar{0}) - B}{\epsilon}\right)$ iterations, ARHT_robust takes $O\left(\log\left(n \log \frac{f(\bar{0}) - B}{\epsilon}\right)\right)$ iterations, and ARHT takes $O\left(\log \frac{f(\bar{0}) - B}{\epsilon}\right)$ iterations. In conclusion, the total number of iterations is $O\left(s \log^2 \frac{f(\bar{0}) - B}{\epsilon} \log\left(n \log \frac{f(\bar{0}) - B}{\epsilon}\right)\right)$, each of which requires a constant number of minimizations of f . \square

Lemma 7.7 (Convergence rate). If Algorithm 1 executes at least $T_1 = s \log \frac{g(x^0) - f(x^*)}{\epsilon}$ Type 1 iterations, then $f(x^T) \leq f(x^*) + \epsilon$.

Proof. By Lemma 7.2, and if we set $\tau = \frac{1}{s}$, in each Type 1 step we have

$$\begin{aligned} g(x^t) - g(x^{t+1}) &\geq \tau (g(x^t) - f(x^*)) \\ \Rightarrow g(x^{t+1}) - f(x^*) &\leq (1 - \tau)(g(x^t) - f(x^*)) \end{aligned}$$

and in each Type 2 step we have

$$g(x^{t+1}) - f(x^*) \leq g(x^t) - f(x^*)$$

(since g can only decrease when unregularizing), therefore

$$\begin{aligned} f(x^T) - f(x^*) &\leq g(x^T) - f(x^*) \\ &\leq (1 - \tau)^{T_1} (g(x^0) - f(x^*)) \\ &\leq e^{-\tau T_1} (g(x^0) - f(x^*)) \\ &\leq \epsilon \end{aligned}$$

where we used the fact that $T_1 = \frac{1}{\tau} \log \frac{g(x^0) - f(x^*)}{\epsilon}$. \square

Proof of Corollary 3.5. By strong convexity we have

$$\begin{aligned} \epsilon &\geq f(x) - f(x^*) \\ &\geq \langle \nabla f(x^*), x - x^* \rangle + \frac{\rho^-}{2} \|x - x^*\|_2^2 \\ &\geq -\zeta \|x - x^*\|_2 + \frac{\rho^-}{2} \|x - x^*\|_2^2 \end{aligned}$$

therefore

$$\frac{\rho^-}{2} \|x - x^*\|_2^2 - \zeta \|x - x^*\|_2 - \epsilon \leq 0$$

looking at which as a quadratic polynomial in $\|x - x^*\|_2$, it follows that

$$\begin{aligned} \|x - x^*\|_2 &\leq \frac{\zeta + \sqrt{\zeta^2 + 2\epsilon\rho^-}}{\rho^-} \\ &= \frac{\zeta}{\rho^-} \left(1 + \sqrt{1 + 2\epsilon\frac{\rho^-}{\zeta^2}} \right) \\ &= (2 + \theta) \frac{\zeta}{\rho^-} \end{aligned}$$

by setting $\epsilon = \frac{\zeta^2}{\rho^-} (\theta + \frac{1}{2}\theta^2)$. \square

Proof of Corollary 3.6. Let us suppose that $S^* \setminus S^t \neq \emptyset$. By restricted strong convexity we have

$$\begin{aligned} &-\frac{1}{2\rho^-}\zeta^2 + \frac{\rho^-}{2}(x_{\min}^*)^2 \\ &> \epsilon \\ &\geq f(x) - f(x^*) \\ &\geq \langle \nabla f(x^*), x - x^* \rangle + \frac{\rho^-}{2} \|x - x^*\|_2^2 \\ &\geq \langle \nabla f(x^*), x \rangle + \frac{\rho^-}{2} \|x_{S^t \setminus S^*}\|_2^2 + \frac{\rho^-}{2} \|x_{S^* \setminus S^t}^*\|_2^2 \\ &\geq -\frac{1}{2\rho^-} \|\nabla_{S^t \setminus S^*} f(x^*)\|_2^2 + \frac{\rho^-}{2} \|x_{S^* \setminus S^t}^*\|_2^2 \\ &\geq -\frac{1}{2\rho^-}\zeta^2 + \frac{\rho^-}{2}(x_{\min}^*)^2 \end{aligned}$$

a contradiction. Here we used the fact that by local optimality $\nabla_{S^*} f(x^*) = \vec{0}$, the inequality $\langle u, v \rangle + \frac{\lambda}{2} \|v\|_2^2 \geq -\frac{1}{2\lambda} \|u\|_2^2$ for any vectors u, v and scalar $\lambda > 0$, and the fact that $\|\nabla_{S^t \setminus S^*} f(x^*)\|_2^2 \leq \zeta^2$ by Definition 2.6. Therefore $S^* \subseteq S^t$. \square

7.2. Analysis of Orthogonal Matching Pursuit with Replacement

We will assume that $s \leq 20\tilde{\kappa}s^*$, as the other case is subsumed by Algorithm 2. Let us denote $\mu = \sqrt{\frac{s^*}{s}}$.

The following technical lemma is at the core of our approach, and roughly states that if there is significant ℓ_2 norm difference between x^t and x^* , at least one of x^t, x^* is significantly larger than \tilde{x}^t in function value. Its importance lies on the fact that instead of directly applying strong convexity between x^t and x^* , it gets a tighter bound by making use of \tilde{x}^t .

Lemma 7.8. For any function f with RSC constant ρ^- at sparsity level $s + s^*$ and any two solutions x^t, x^* with respective supports S^t, S^* and sparsity levels s, s^* , we have

that

$$\begin{aligned} &(\sqrt{f(x^t) - f(\tilde{x}^t)} + \sqrt{f(x^*) - f(\tilde{x}^t)})^2 \\ &\geq \frac{\rho^-}{2} \left(\|x_{S^* \setminus S^t}^*\|_2^2 + \|x_{S^t \setminus S^*}^t\|_2^2 \right) \end{aligned}$$

Proof. We have

$$\begin{aligned} &(\sqrt{f(x^t) - f(\tilde{x}^t)} + \sqrt{f(x^*) - f(\tilde{x}^t)})^2 \\ &\geq \frac{\rho^-}{2} (\|x^t - \tilde{x}^t\|_2 + \|x^* - \tilde{x}^t\|_2)^2 \\ &\geq \frac{\rho^-}{2} \|x^t - x^*\|_2^2 \\ &\geq \frac{\rho^-}{2} \left(\|x_{S^* \setminus S^t}^*\|_2^2 + \|x_{S^t \setminus S^*}^t\|_2^2 \right) \end{aligned}$$

where the first inequality follows by applying strong convexity to lower bound $f(x^t) - f(\tilde{x}^t)$ and $f(x^*) - f(\tilde{x}^t)$ combined with the fact that by definition of \tilde{x}^t , $\nabla_{S^t \cup S^*} f(\tilde{x}^t) = \vec{0}$, and the second is a triangle inequality. \square

We now provide the main technical lemma for this section, which bounds the progress of Algorithm 2 in one iteration.

Lemma 7.9. We can bound the progress of one step of the algorithm by distinguishing the following three cases:

• If $\mu\tilde{\kappa} \leq 1$, then

$$f(x^{t+1}) - f(x^*) \leq (f(x^t) - f(x^*)) \left(1 - \frac{\mu}{|S^* \setminus S^t|} \right)$$

• If $\mu\tilde{\kappa} > 1$ and $f(x^*) = f(\tilde{x}^t)$, then

$$\begin{aligned} &f(x^{t+1}) - f(x^*) \leq (f(x^t) - f(x^*)) \\ &\cdot \left(1 - \frac{\mu}{|S^* \setminus S^t|} (2 - \mu\tilde{\kappa}) \right) \end{aligned}$$

• If $\mu\tilde{\kappa} > 1$ and $f(x^*) > f(\tilde{x}^t)$, then

$$\begin{aligned} &f(x^{t+1}) - f(x^*) \leq (f(x^t) - f(x^*)) \\ &\cdot \left(1 - \frac{\mu}{|S^* \setminus S^t|} \left(2 - \mu\tilde{\kappa} - \frac{2(\mu\tilde{\kappa} - 1)}{\sqrt{\frac{f(x^t) - f(\tilde{x}^t)}{f(x^*) - f(\tilde{x}^t)} - 1}} \right) \right) \end{aligned}$$

Proof. First of all, if $S^* \subseteq S^t$ then, since x^t is an S^t -restricted minimizer, we have $f(x^t) \leq f(x^*) \leq f(x^*)$ and we are done. So suppose otherwise, i.e. $S^* \setminus S^t \neq \emptyset$ and $f(x^t) > f(x^*)$. Let $i = \operatorname{argmax}_{i \notin S^t} |\nabla_i f(x^t)|$ and $j =$

$\operatorname{argmin}_{j \in S^t} |x_j^t|$. By definition of OMPR (Algorithm 3) and restricted smoothness of f , we have

$$\begin{aligned}
 & f(x^{t+1}) \\
 & \leq \min_{\eta \in \mathbb{R}} f(x^t + \eta \vec{1}_i - x_j^t \vec{1}_j) \\
 & \leq \min_{\eta \in \mathbb{R}} f(x^t) + \langle \nabla f(x^t), \eta \vec{1}_i - x_j^t \vec{1}_j \rangle + \frac{\rho^+}{2} \left\| \eta \vec{1}_i - x_j^t \vec{1}_j \right\|_2^2 \\
 & = \min_{\eta \in \mathbb{R}} f(x^t) + \eta \nabla_i f(x^t) + \frac{\rho^+}{2} \eta^2 + \frac{\rho^+}{2} (x_j^t)^2 \\
 & = f(x^t) - \frac{(\nabla_i f(x^t))^2}{2\rho^+} + \frac{\rho^+}{2} (x_j^t)^2 \\
 & \leq f(x^t) - \frac{\|\nabla_{S^* \setminus S^t} f(x^t)\|_2^2}{2\rho^+ |S^* \setminus S^t|} + \frac{\rho^+}{2|S^t \setminus S^*|} \left\| x_{S^t \setminus S^*}^t \right\|_2^2
 \end{aligned}$$

where the second to last equality follows from the fact that $\nabla_j f(x^t) = \vec{0}$, as x^t is an S^t -restricted minimizer of f , and the last inequality since

$$(x_j^t)^2 = \min_{j \in S^t \setminus S^*} (x_j^t)^2 \leq \frac{\left\| x_{S^t \setminus S^*}^t \right\|_2^2}{|S^t \setminus S^*|}$$

Re-arranging, we get

$$\begin{aligned}
 & |S^* \setminus S^t| (f(x^t) - f(x^{t+1})) \\
 & \geq \frac{\|\nabla_{S^* \setminus S^t} f(x^t)\|_2^2}{2\rho^+} - \frac{\rho^+}{2} \frac{|S^* \setminus S^t|}{|S^t \setminus S^*|} \left\| x_{S^t \setminus S^*}^t \right\|_2^2 \quad (9)
 \end{aligned}$$

On the other hand, by restricted strong convexity of f ,

$$\begin{aligned}
 & f(x^*) - f(x^t) \\
 & \geq \langle \nabla f(x^t), x^* - x^t \rangle + \frac{\rho^-}{2} \|x^* - x^t\|_2^2 \\
 & = \langle \nabla_{S^* \setminus S^t} f(x^t), x_{S^* \setminus S^t}^* \rangle + \frac{\rho^-}{2} \|x^* - x^t\|_2^2 \\
 & \geq \langle \nabla_{S^* \setminus S^t} f(x^t), x_{S^* \setminus S^t}^* \rangle + \frac{\rho^-}{2} \left\| x_{S^* \setminus S^t}^* \right\|_2^2 + \frac{\rho^-}{2} \left\| x_{S^t \setminus S^*}^t \right\|_2^2 \\
 & \geq \langle \nabla_{S^* \setminus S^t} f(x^t), x_{S^* \setminus S^t}^* \rangle + \frac{\mu\rho^+}{2} \left\| x_{S^* \setminus S^t}^* \right\|_2^2 \\
 & + \frac{\rho^- - \mu\rho^+}{2} \left\| x_{S^* \setminus S^t}^* \right\|_2^2 + \frac{\rho^-}{2} \left\| x_{S^t \setminus S^*}^t \right\|_2^2 \\
 & \geq -\frac{1}{2\mu\rho^+} \left\| \nabla_{S^* \setminus S^t} f(x^t) \right\|_2^2 \\
 & + \frac{\rho^- - \mu\rho^+}{2} \left\| x_{S^* \setminus S^t}^* \right\|_2^2 + \frac{\rho^-}{2} \left\| x_{S^t \setminus S^*}^t \right\|_2^2
 \end{aligned}$$

where the first equality follows from the fact that $\nabla_{S^t} f(x^t) = \vec{0}$ as x^t is an S^t -restricted minimizer of f and the last inequality from using the fact that $\langle u, v \rangle + \frac{\lambda}{2} \|v\|_2^2 \geq -\frac{1}{2\lambda} \|u\|_2^2$ for any $\lambda > 0$.

Re-arranging, we get

$$\begin{aligned}
 & \frac{1}{2\mu\rho^+} \left\| \nabla_{S^* \setminus S^t} f(x^t) \right\|_2^2 \\
 & \geq f(x^t) - f(x^*) - \frac{\mu\rho^+ - \rho^-}{2} \left\| x_{S^* \setminus S^t}^* \right\|_2^2 + \frac{\rho^-}{2} \left\| x_{S^t \setminus S^*}^t \right\|_2^2
 \end{aligned}$$

By substituting this into (10),

$$\begin{aligned}
 & |S^* \setminus S^t| (f(x^t) - f(x^{t+1})) \\
 & \geq \mu (f(x^t) - f(x^*)) - \frac{\mu^2\rho^+ - \mu\rho^-}{2} \left\| x_{S^* \setminus S^t}^* \right\|_2^2 + \frac{\mu\rho^-}{2} \left\| x_{S^t \setminus S^*}^t \right\|_2^2 \\
 & - \frac{\rho^+}{2} \frac{|S^* \setminus S^t|}{|S^t \setminus S^*|} \left\| x_{S^t \setminus S^*}^t \right\|_2^2
 \end{aligned}$$

Note that by our choice of μ and since $s^* \leq s$,

$$\mu^2\rho^+ = \rho^+ \frac{s^*}{s} \geq \rho^+ \frac{s^* - |S^* \cap S^t|}{s - |S^* \cap S^t|} = \rho^+ \frac{|S^* \setminus S^t|}{|S^t \setminus S^*|}$$

and so

$$\begin{aligned}
 & \mu (f(x^t) - f(x^*)) - \frac{\mu^2\rho^+ - \mu\rho^-}{2} \left\| x_{S^* \setminus S^t}^* \right\|_2^2 + \frac{\mu\rho^-}{2} \left\| x_{S^t \setminus S^*}^t \right\|_2^2 \\
 & - \frac{\rho^+}{2} \frac{|S^* \setminus S^t|}{|S^t \setminus S^*|} \left\| x_{S^t \setminus S^*}^t \right\|_2^2 \\
 & \geq \mu (f(x^t) - f(x^*)) \\
 & - \frac{\mu}{2} (\mu\rho^+ - \rho^-) \left(\left\| x_{S^* \setminus S^t}^* \right\|_2^2 + \left\| x_{S^t \setminus S^*}^t \right\|_2^2 \right)
 \end{aligned}$$

concluding that

$$\begin{aligned}
 & |S^* \setminus S^t| (f(x^t) - f(x^{t+1})) \\
 & \geq \mu (f(x^t) - f(x^*)) \\
 & - \frac{\mu}{2} (\mu\rho^+ - \rho^-) \left(\left\| x_{S^* \setminus S^t}^* \right\|_2^2 + \left\| x_{S^t \setminus S^*}^t \right\|_2^2 \right)
 \end{aligned}$$

For $\mu\tilde{\kappa} \leq 1 \Leftrightarrow \mu\rho^+ - \rho^- \leq 0$, this automatically implies that

$$\begin{aligned}
 & f(x^t) - f(x^{t+1}) \geq \frac{\mu}{|S^* \setminus S^t|} (f(x^t) - f(x^*)) \\
 & \Leftrightarrow f(x^{t+1}) - f(x^*) \leq \left(1 - \frac{\mu}{|S^* \setminus S^t|} \right) (f(x^t) - f(x^*))
 \end{aligned}$$

On the other hand, if $\mu\tilde{\kappa} > 1$ we have

$$\begin{aligned}
 & |S^* \setminus S^t| (f(x^t) - f(x^{t+1})) \\
 & \geq \mu (f(x^t) - f(x^*)) \\
 & - \frac{\mu}{2} (\mu\rho^+ - \rho^-) \left(\left\| x_{S^* \setminus S^t}^* \right\|_2^2 + \left\| x_{S^t \setminus S^*}^t \right\|_2^2 \right) \\
 & \geq \mu (f(x^t) - f(x^*)) \\
 & - \mu (\mu\tilde{\kappa} - 1) \left(\sqrt{f(x^t) - f(\tilde{x}^t)} + \sqrt{f(x^*) - f(\tilde{x}^t)} \right)^2
 \end{aligned}$$

where we used Lemma 7.8. If $f(x^*) = f(\tilde{x}^t)$ it is immediate that

$$f(x^{t+1}) - f(x^*) \leq \left(1 - \frac{\mu}{|S^* \setminus S^t|} (2 - \mu\tilde{\kappa})\right) (f(x^t) - f(x^*))$$

so let us from now on assume that $f(x^*) > f(\tilde{x}^t)$ and set $a = f(x^t) - f(\tilde{x}^t)$, $a' = f(x^{t+1}) - f(\tilde{x}^t)$, and $b = f(x^*) - f(\tilde{x}^t)$. From what we have concluded before

$$|S^* \setminus S^t| (a - a') \geq \mu(a - b) - \mu(\mu\tilde{\kappa} - 1) (\sqrt{a} + \sqrt{b})^2$$

or equivalently

$$\begin{aligned} a' - b &\leq \left(1 - \frac{\mu}{|S^* \setminus S^t|}\right) (a - b) + \frac{\mu}{|S^* \setminus S^t|} (\mu\tilde{\kappa} - 1) (\sqrt{a} + \sqrt{b})^2 \\ &= (a - b) \left(1 - \frac{\mu}{|S^* \setminus S^t|} \left(1 - (\mu\tilde{\kappa} - 1) \frac{(\sqrt{a} + \sqrt{b})^2}{a - b}\right)\right) \\ &= (a - b) \left(1 - \frac{\mu}{|S^* \setminus S^t|} \left(1 - (\mu\tilde{\kappa} - 1) \frac{\sqrt{\frac{a}{b}} + 1}{\sqrt{\frac{a}{b}} - 1}\right)\right) \\ &= (a - b) \left(1 - \frac{\mu}{|S^* \setminus S^t|} \left(1 - (\mu\tilde{\kappa} - 1) \left(1 + \frac{2}{\sqrt{\frac{a}{b}} - 1}\right)\right)\right) \\ &= (a - b) \left(1 - \frac{\mu}{|S^* \setminus S^t|} \left(2 - \mu\tilde{\kappa} - \frac{2(\mu\tilde{\kappa} - 1)}{\sqrt{\frac{a}{b}} - 1}\right)\right) \end{aligned}$$

Replacing back a, a', b , the desired statement follows:

$$\begin{aligned} f(x^{t+1}) - f(x^*) &\leq (f(x^t) - f(x^*)) \\ &\cdot \left(1 - \frac{\mu}{|S^* \setminus S^t|} \left(2 - \mu\tilde{\kappa} - \frac{2(\mu\tilde{\kappa} - 1)}{\sqrt{\frac{f(x^t) - f(\tilde{x}^t)}{f(x^*) - f(\tilde{x}^t)}} - 1}\right)\right) \end{aligned}$$

□

Having established the above lemma, the proof of Theorem 3.7 follows easily.

Proof of Theorem 3.7.

Case 1: $\mu\tilde{\kappa} \leq 1$. By Lemma 7.9, we have

$$\begin{aligned} f(x^T) - f(x^*) &\leq (f(x^{T-1}) - f(x^*)) \left(1 - \frac{\mu}{|S^* \setminus S^{T-1}|}\right) \\ &\leq (f(x^{T-1}) - f(x^*)) \left(1 - \frac{\mu}{s^*}\right) \\ &\leq (f(x^{T-1}) - f(x^*)) e^{-\frac{\mu}{s^*}} \\ &\leq \dots \\ &\leq (f(x^0) - f(x^*)) e^{-T \frac{\mu}{s^*}} \\ &\leq \epsilon \end{aligned}$$

for our choice of $T = O\left(\sqrt{ss^*} \log \frac{f(x^0) - f(x^*)}{\epsilon}\right)$ and replacing $\mu = \sqrt{\frac{s^*}{s}}$.

Case 2: $\mu\tilde{\kappa} > 1$. Let \mathcal{A} be the set of $0 \leq t \leq T - 1$ such that $f(x^*) = f(\tilde{x}^t)$ and \mathcal{B} the set of $0 \leq t \leq T - 1$ such that $f(x^*) > f(\tilde{x}^t)$. By Lemma 7.9, for $t \in \mathcal{A}$ we then have

$$\begin{aligned} f(x^{t+1}) - f(x^*) &\leq (f(x^t) - f(x^*)) \left(1 - \frac{\mu}{|S^* \setminus S^t|} (2 - \mu\tilde{\kappa})\right) \\ &\leq (f(x^t) - f(x^*)) \left(1 - \frac{\mu}{s^*} (2 - \mu\tilde{\kappa})\right) \end{aligned}$$

We now consider the case $t \in \mathcal{B}$. By Lemma 7.9,

$$\begin{aligned} f(x^{t+1}) - f(x^*) &\leq (f(x^t) - f(x^*)) \\ &\cdot \left(1 - \frac{\mu}{|S^* \setminus S^t|} \left(2 - \mu\tilde{\kappa} - \frac{2(\mu\tilde{\kappa} - 1)}{\sqrt{\frac{f(x^t) - f(\tilde{x}^t)}{f(x^*) - f(\tilde{x}^t)}} - 1}\right)\right) \end{aligned} \tag{11}$$

Let us suppose that the Theorem statement is not true. This implies

$$\begin{aligned} f(x^t) - f(x^*) &\geq f(x^T) - f(x^*) \\ &> \epsilon + \frac{4(1 - \theta)(\mu\tilde{\kappa} - 1)}{(2 - \mu\tilde{\kappa} - \theta)^2} (f(x^*) - f(\tilde{x}^*)) \\ &\geq \epsilon + \frac{4(1 - \theta)(\mu\tilde{\kappa} - 1)}{(2 - \mu\tilde{\kappa} - \theta)^2} (f(x^*) - f(\tilde{x}^t)) \\ &\geq \frac{4(1 - \theta)(\mu\tilde{\kappa} - 1)}{(2 - \mu\tilde{\kappa} - \theta)^2} (f(x^*) - f(\tilde{x}^t)) \end{aligned} \tag{12}$$

for all $0 \leq t \leq T$. Therefore

$$\begin{aligned}
 & f(x^t) - f(\tilde{x}^t) \\
 & > \left(\frac{4(1-\theta)(\mu\tilde{\kappa}-1)}{(2-\mu\tilde{\kappa}-\theta)^2} + 1 \right) (f(x^*) - f(\tilde{x}^t)) \\
 & = \left(\frac{4(1-\theta)(\mu\tilde{\kappa}-1) + 4 + (\mu\tilde{\kappa} + \theta)^2 - 4(\mu\tilde{\kappa} + \theta)}{(2-\mu\tilde{\kappa}-\theta)^2} \right) \\
 & \cdot (f(x^*) - f(\tilde{x}^t)) \\
 & = \frac{(\mu\tilde{\kappa}-\theta)^2}{(2-\mu\tilde{\kappa}-\theta)^2} (f(x^*) - f(\tilde{x}^t))
 \end{aligned}$$

or equivalently for all $t \in \mathcal{B}$

$$\sqrt{\frac{f(x^t) - f(\tilde{x}^t)}{f(x^*) - f(\tilde{x}^t)}} - 1 > \frac{\mu\tilde{\kappa} - \theta}{2 - \mu\tilde{\kappa} - \theta} - 1 = \frac{2(\mu\tilde{\kappa} - 1)}{2 - \mu\tilde{\kappa} - \theta}$$

Replacing this into (11), we get that for any $t \in \mathcal{B}$

$$\begin{aligned}
 & f(x^{t+1}) - f(x^*) \\
 & \leq (f(x^t) - f(x^*)) \\
 & \cdot \left(1 - \frac{\mu}{|S^* \setminus S^t|} \left(2 - \mu\tilde{\kappa} - \frac{2(\mu\tilde{\kappa}-1)}{\sqrt{\frac{f(x^t)-f(\tilde{x}^t)}{f(x^*)-f(\tilde{x}^t)}} - 1}} \right) \right) \\
 & \leq (f(x^t) - f(x^*)) \left(1 - \frac{\mu}{|S^* \setminus S^t|} \theta \right)
 \end{aligned}$$

and so combining it with the case $t \in \mathcal{A}$ and using the fact that $\mu\tilde{\kappa} < 2 - \theta \Leftrightarrow \theta < 2 - \mu\tilde{\kappa}$,

$$\begin{aligned}
 & f(x^T) - f(x^*) \\
 & \leq (f(x^{T-1}) - f(x^*)) \left(1 - \frac{\mu}{|S^* \setminus S^{T-1}|} \min\{2 - \mu\tilde{\kappa}, \theta\} \right) \\
 & \leq (f(x^{T-1}) - f(x^*)) \left(1 - \frac{\mu}{s^*} \theta \right) \\
 & \leq (f(x^{T-1}) - f(x^*)) e^{-\frac{\mu}{s^*} \theta} \\
 & \leq \dots \\
 & \leq (f(x^0) - f(x^*)) e^{-T \frac{\mu}{s^*} \theta} \\
 & = \epsilon + \frac{4(1-\theta)(\mu\tilde{\kappa}-1)}{(2-\mu\tilde{\kappa}-\theta)^2} (f(x^*) - f(\tilde{x}^*))
 \end{aligned}$$

where the last equality follows by our choice of

$$T = \frac{\sqrt{ss^*}}{\theta} \log \frac{f(x^0) - f(x^*)}{B}$$

and replacing $\mu = \sqrt{\frac{s^*}{s}}$. This is a contradiction. \square

8. Experiments

For experimental evaluation we used well known and publicly available datasets. Their names and basic properties are outlined in Table 2.

Table 2. Datasets used for experimental evaluation. The columns are the dataset name, the number of examples m , and the number of features n . The datasets can be downloaded [here](#).

NAME	n	d	PROBLEM
KDDCUP04_BIO	145750	74	BINARY
CAL_HOUSING	20639	8	REGRESSION
CENSUS	299284	401	BINARY
COMP-ACTIV-HARDER	8191	12	REGRESSION
IJCNN1	24995	22	BINARY
LETTER	20000	16	BINARY
SLICE	53500	384	REGRESSION
YEAR	463715	90	REGRESSION

8.1. Setup details

8.1.1. BASIC DEFINITIONS

The two quantities that take part in our experiments are the *sparsity* and the *loss* of a particular solution. We have already defined and discussed the former at length. The latter refers to the training loss for the problems of Linear Regression and Logistic Regression. We let m denote the number of examples and n the number of features in each example.

In the *Linear Regression* task we are given the dataset (A, b) , where $A \in \mathbb{R}^{m \times n}$, $b \in \mathbb{R}^m$. The columns of A correspond to features and the rows to examples. The (ℓ_2 *Linear Regression*) loss of a solution $x \in \mathbb{R}^n$ is defined as $\ell_2\text{-loss}(x) = \frac{1}{2} \|Ax - b\|_2^2$.

In the *Logistic Regression* task we are given the dataset (A, b) , where $A \in \mathbb{R}^{m \times n}$, $b \in \{0, 1\}^m$. The columns of A correspond to features and the rows to examples. The (*Logistic Regression*) loss of a solution $x \in \mathbb{R}^n$ is defined as $\text{logistic_loss}(x) = \sum_{i \in [m]} (-b_i \log \sigma(Ax)_i - (1 - b_i) \log(1 - \sigma(Ax)_i))$, where $\sigma : \mathbb{R} \rightarrow \mathbb{R}$ defined as $\sigma(t) = \frac{1}{1 + e^{-t}}$ is the sigmoid function.

8.1.2. DATA PRE-PROCESSING

We apply a very basic form of pre-processing to the data. More specifically, we use one-hot encoding to turn categorical features into numerical ones. Then, we discard any examples with missing data so that all the entries of A are defined. We also augment the matrix A with an extra all-ones column (i.e. $\vec{1}$) in order to encode the constant (y -intercept) term into A , and we scale all the columns of A so that their ℓ_2 norm is 1. Finally, for the case of ARHT we further augment A in order to encode the regularizer as well. We do this by adding an identity matrix as extra rows. In other words, $A \leftarrow \begin{pmatrix} A \\ I \end{pmatrix}$ and $b \leftarrow \begin{pmatrix} b \\ \vec{0} \end{pmatrix}$.

8.2. Implementation details

The code has been implemented in *python3*, with libraries *numpy*, *sklearn*, and *scipy*.

8.2.1. INNER OPTIMIZATION PROBLEM

All the algorithms except for LASSO rely on an inner optimization routine in a restricted subset of coordinates in each step. The inner optimization problem consists of solving a standard Linear Regression or Logistic Regression problem using only a submatrix of A defined by a subset of s of its columns. For that, we use *LinearRegression* and *LogisticRegression* from *sklearn.linear_model*. For Logistic Regression we used an LBFGS solver with 1000 iterations.

8.2.2. OVERALL ALGORITHM

The LASSO solver we used is *Lasso* from *sklearn.linear_model* with 1000 iterations. As LASSO is not tuned in terms of a required sparsity s , but rather in terms of the regularization parameter α , for each sparsity level we applied binary search on α in order to find a parameter α that gives the required sparsity.

For ARHT, we used a fixed number of 20 iterations at Line 5 of Algorithm 2. In Line 19 of Algorithm 1 we slightly weaken the progress condition to

$$g_{R^t}(x^t) - g_{R^t}(x^{t+1}) \geq \frac{10^{-3}}{s} (g_{R^t}(x^t) - \text{opt}) \quad (13)$$

so that it does not depend Furthermore, we do not perform a fixed number of iterations. Instead, we use a stopping criterion: If the progress condition (13) is not met and at least half the elements in x^t have already been unregularized, i.e. $|S^t \setminus R^t| \geq \frac{1}{2} |S^t|$, then we stop. If a desirable solution has not been found, it means that this might be an unsuccessful run, and early termination can be used to detect such runs early and re-start, thus improving the runtime. The routine which samples an index i proportional to x_i^2 was implementing by a standard sampling method that uses binary search on i and flips a random coin at each step. This requires computation of interval sums of x_i^2 , which is done by computing partial sums.

8.3. Additional experiments

In Figure 3 and Figure 4 we present some additional experiments, as an addendum to Section 4. One can observe that in these datasets the performance of all algorithms is very close, except for LASSO which performs worse in some cases.

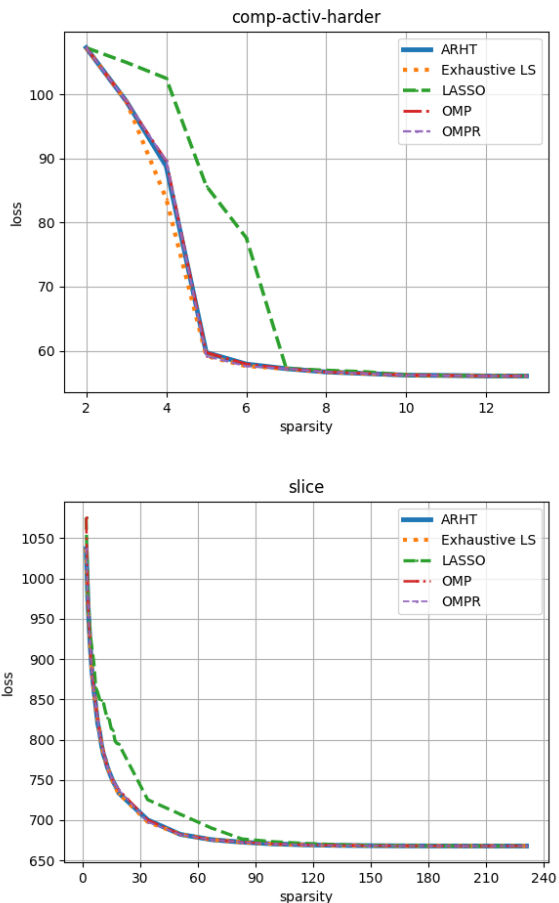


Figure 3. Comparison of different algorithms in the Regression datasets *comp-activ-harder* and *slice* using the Linear Regression loss.

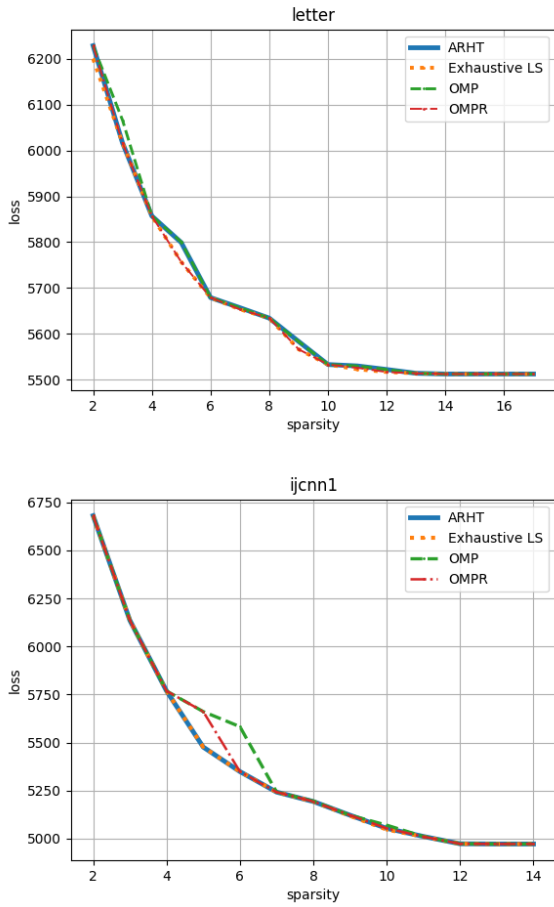


Figure 4. Comparison of different algorithms in the Binary classification datasets *letter* and *ijcnn1* using the Logistic Regression loss.