

# GRADIENT DESCENT CONVERGES LINEARLY FOR LOGISTIC REGRESSION ON SEPARABLE DATA

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

We show that running gradient descent on the logistic regression objective guarantees loss  $f(\mathbf{x}) \leq 1.1 \cdot f(\mathbf{x}^*) + \varepsilon$ , where the error  $\varepsilon$  decays exponentially with the number of iterations. This is in contrast to the common intuition that the absence of strong convexity precludes linear convergence of first-order methods, and highlights the importance of variable learning rates for gradient descent. For separable data, our analysis proves that the error between the predictor returned by gradient descent and the hard SVM predictor decays as  $\text{poly}(1/t)$ , exponentially faster than the previously known bound of  $O(\log \log t / \log t)$ . Our key observation is a property of the logistic loss that we call multiplicative smoothness and is (surprisingly) little-explored: As the loss decreases, the objective becomes (locally) smoother and therefore the learning rate can increase. Our results also extend to sparse logistic regression, where they lead to an exponential improvement of the sparsity-error tradeoff.

## 1 INTRODUCTION

Logistic regression is one of the most widely used classification methods because of its simplicity, interpretability, and good practical performance. Yet, the convergence behavior of first-order methods on this task is not well understood: In practice gradient descent performs much better than what the theory predicts. In particular, a general analysis of gradient descent for smooth functions implies convergence with the error in function value decaying as  $O(1/T)$ . Analyses with stronger, linear convergence guarantees generally require the function to satisfy the strong convexity property, which, in contrast to other losses such as the  $\ell_2$  loss, the logistic loss only satisfies in a bounded set of solutions around zero. As a result, this introduces an *exponential* runtime dependency on the magnitude of the optimal solution Rätsch et al. (2001); Freund et al. (2018), which is undesirable in practice. This poses a serious obstacle to obtaining favorable error rates for logistic regression that lead to high-precision solutions.

A deeper study into the structure of the exponential and logistic losses was done in Telgarsky & Singer (2012), who showed that, for linearly separable data, greedy coordinate descent achieves linear convergence with a rate that depends on the maximum linear classification margin (i.e. hard SVM margin). Unfortunately, for logistic regression, it also has a  $2^m$  dependence on the number of examples, making it inefficient for any real-world task. The significance of the separability of the data for convergence has also been observed in Telgarsky (2013); Freund et al. (2018), who present convergence results based on quantitative measures of separability. Telgarsky (2013) also refines the results of Telgarsky & Singer (2012) for the exponential loss, but still suffers from an exponential overhead originating the multiplicative discrepancy between the exponential and the logistic loss. Interestingly Telgarsky (2013) points out that logistic regression experiments paint a much more favorable picture than the theory predicts. For separable data, Soudry et al. (2018) showed that the gradient descent logistic regression estimator converges to the maximum margin estimator at a rate of  $O(\log \log T / \log T)$ , which implies function value convergence at a rate of  $O(1/T)$ . Interestingly, Nacson et al. (2019) experimentally observed that these rates seem to be exponentially improvable if one uses variable step sizes, in the case of logistic regression and shallow neural networks. However, as shown in Ji & Telgarsky (2018), the separability assumption is important, and the  $\text{poly}(1/T)$  bound of function value convergence is tight for gradient descent on arbitrary data.

Another approach to obtain high-precision solutions is by using second order methods, which in addition to first order (gradient) information, use second order (Hessian) information about the function. These make use of second order stability properties, such as quasi-self-concordance Bach (2010) combined with Newton’s method Karimireddy et al. (2018), or ball oracles Carmon et al. (2020); Adil et al. (2021). Such approaches are generally not suitable for large-scale applications because of their reliance on repeated calls to large linear system solvers.

**Our work.** In this paper, we show that (under appropriate assumptions) we can get the best of both worlds of first and second order methods, thus giving a partial explanation for the excellent performance that first-order methods have on logistic regression in practice. In particular, given a binary classification instance ( $\mathbf{A} \in \{-1, 1\}^{m \times n}$ ,  $\mathbf{b} \in \{-1, 1\}^m$ ) with associated logistic loss  $f(\mathbf{x}) = \sum_i \log(1 + \exp(-b_i(\mathbf{A}\mathbf{x})_i))$ , we show that simple variants of gradient descent

return a solution with  $f(\mathbf{x}) \leq (1 + \delta) \cdot f(\mathbf{x}^*) + \varepsilon$  after  $O\left(K\left(\frac{1}{\delta} + \log \frac{f(\mathbf{0})}{\varepsilon}\right)\right)$  iterations, where  $K = \text{poly}(n, \|\mathbf{x}^*\|)$ . Even though the error still decays as  $1/T$  in the worst case because of the  $\frac{1}{\delta}$  dependence, the additive error is now  $\delta f(\mathbf{x}^*)$  instead of  $\delta f(\mathbf{0})$ , allowing for much faster convergence when the optimal loss  $f(\mathbf{x}^*)$  is smaller (which is our measure of linear separability of the data). For linearly separable data, i.e. as  $f(\mathbf{x}^*)$  approaches 0, the convergence becomes linear. We also show that the distance to the maximum margin estimator  $\left\| \frac{\mathbf{x}}{\|\mathbf{x}\|_2} - \frac{\mathbf{x}^*}{\|\mathbf{x}^*\|_2} \right\|_2$  decays as  $1/T$ , exponentially improving over the  $\log \log T / \log T$  bound of Soudry et al. (2018).

Instead of properties like Lipschitzness, smoothness, strong convexity that are commonly used in the study of first order methods, we find that there are two properties that are more relevant to the structure of the logistic regression problem. The first one is *second order robustness*, which means that the Hessian is stable (in a spectral sense) in any small enough norm ball Cohen et al. (2017). This is closely related to quasi-self-concordance, a property that has been previously used in the analysis of second order algorithms Bach (2010). The second property is what we call *multiplicative smoothness*, which means that the function is locally smooth, with the smoothness constant being proportional to the function value (loss). Together, these properties show that, as the loss decreases, the objective becomes (locally) smoother and therefore the learning rate can increase. This motivates a variable step size schedule that is inversely proportional to the loss, thus making larger steps as the solution approaches optimality. This in fact agrees with the observations of Soudry et al. (2018); Nacson et al. (2019) on the importance of a variable learning rate. As can be seen in the toy example from Soudry et al. (2018) in Figure 1, simply replacing the fixed learning rate  $\eta$  used in Soudry et al. (2018) by an increasing learning rate  $\eta \cdot f(\mathbf{x}^0)/f(\mathbf{x}^t)$  yields an exponential improvement, both in loss and distance to the maximum margin estimator.

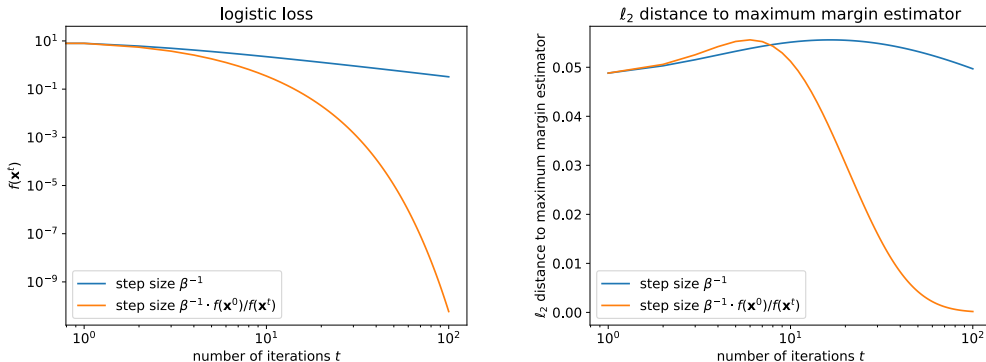


Figure 1: Comparison between fixed and increasing step sizes in the toy example from Figure 1 of Soudry et al. (2018). The fixed step size is set to  $\beta^{-1} := \|\mathbf{A}\|_2^{-2}$ , and the increasing to  $\beta^{-1} f(\mathbf{x}^0)/f(\mathbf{x}^t)$ . The estimator error is defined as  $\|\mathbf{x}^t / \|\mathbf{x}^t\|_2 - \mathbf{x}^* / \|\mathbf{x}^*\|_2\|_2$ .

Algorithm	Order	Guarantee	Runtime Error Dependence
Gradient descent	First	$f(\mathbf{x}) \leq f(\mathbf{x}^*) + \varepsilon$	$m/\varepsilon$
Accelerated gradient descent	First	$f(\mathbf{x}) \leq f(\mathbf{x}^*) + \varepsilon$	$\sqrt{m/\varepsilon}$
Newton/Trust region	Second	$f(\mathbf{x}) \leq f(\mathbf{x}^*) + \varepsilon$	$\log(m/\varepsilon)$
This paper	First	$f(\mathbf{x}) \leq (1 + \delta) \cdot f(\mathbf{x}^*) + \varepsilon$	$\delta^{-1} + \log(m/\varepsilon)$

Table 1: Algorithms for logistic regression and dependence on  $m/\varepsilon$  (omitting extra  $\text{polylog}(m, n)$  factors). Algorithms with exponential dependences on any problem parameter are omitted.

Algorithm	Guarantee	Sparsity	Order
Shalev-Shwartz et al. (2010)	$f(\mathbf{x}) \leq f(\mathbf{x}^*) + \varepsilon$	$\ \mathbf{x}^*\ _1^2 m/\varepsilon$	First
This paper	$f(\mathbf{x}) \leq (1 + \delta) \cdot f(\mathbf{x}^*) + \varepsilon$	$\ \mathbf{x}^*\ _1^2 (\delta^{-1} + \log(m/\varepsilon))$	First

Table 2: Algorithms for sparse logistic regression

### 1.1 SPARSE LOGISTIC REGRESSION

In practice, it is often important to force the solution of a logistic regression problem to be *sparse*, i.e. have only a few non-zero entries, which is a form of feature selection. This is because most of the features might only be marginally useful, and thus one can drastically reduce the size of the model while not significantly sacrificing the predictive performance. Apart from computational efficiency, feature selection is also important to improve interpretability and avoid overfitting.

Most progress in sparse optimization has focused on objective functions with condition number bounded by some  $\kappa > 0$ . Results in this line of work guarantee a solution with relaxed sparsity  $s' \geq s$ , where  $s$  is the target sparsity, and algorithms include lasso, orthogonal matching pursuit (OMP), and iterative hard thresholding (IHT) Natarajan (1995); Blumensath & Davies (2009); Shalev-Shwartz et al. (2010); Jain et al. (2011; 2014); Axiotis & Sviridenko (2021; 2022). The state of the art result by Axiotis & Sviridenko (2022) gives a sparsity of  $s' = O(\kappa) \cdot s$  using a variant of the IHT algorithm.

However, the condition number of the logistic loss is unbounded, because it is not strongly convex. Therefore, these results do not directly apply, although they do apply to  $\ell_2$ -regularized logistic regression. Some works Van de Geer (2008); Bunea (2008) have analyzed lasso methods for logistic regression without condition number assumptions, and Shalev-Shwartz et al. (2010) provides three different analyses for smooth but not strongly convex functions. These apply to logistic regression and give a sparsity of  $O\left(\|\mathbf{x}^*\|_1^2 \frac{m}{\varepsilon}\right)$  to achieve a loss of  $f(\mathbf{x}) \leq f(\mathbf{x}^*) + \varepsilon$ . The most practical of these is a forward greedy selection algorithm, which is also known as greedy coordinate descent.

**Our work.** Using the second order stability and multiplicative smoothness properties, we show that a slight variation of greedy coordinate descent gives a sparsity of

$$O\left(\|\mathbf{x}^*\|_1^2 (\delta^{-1} + \log(m/\varepsilon))\right)$$

and a loss of  $f(\mathbf{x}) \leq (1 + \delta) \cdot f(\mathbf{x}^*) + \varepsilon$ . As long as the  $1 + \delta$  approximation in front of  $f(\mathbf{x}^*)$  is tolerated, as is the case when  $f(\mathbf{x}^*) \ll m$ , this implies an exponential improvement in the  $\varepsilon$  dependence from  $\frac{m}{\varepsilon}$  to  $\log \frac{m}{\varepsilon}$ . In addition, our analysis does not require (but is also not affected by) fully corrective steps, in which the function is fully re-optimized over the support of the current solution.

## 2 PRELIMINARIES

**Notation.** We denote  $[n] = \{1, 2, \dots, n\}$ . We will use **bold** to refer to vectors or matrices. We denote by  $\mathbf{0}$  the all-zero vector,  $\mathbf{1}$  the all-one vector,  $\mathbf{O}$  the all-zero matrix, and by  $\mathbf{I}$  the identity

matrix (with dimensions understood from the context). Additionally, we will denote by  $\mathbf{1}_i$  the  $i$ -th basis vector, i.e. the vector that is 0 everywhere except at position  $i$ .

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

For any vector  $\mathbf{x} \in \mathbb{R}^n$  and set  $S \subseteq [n]$ , we denote by  $\mathbf{x}_S$  the vector that results from  $\mathbf{x}$  after zeroing out all the entries except those in positions given by indices in  $S$ . We will also use the notation  $\nabla_S f(\mathbf{x}) := (\nabla f(\mathbf{x}))_S$  to denote the restriction of a gradient to  $S$ .

We use the notation  $\tilde{O}(\cdot)$  to hide  $\text{poly} \log(n, m)$  factors in  $O$ -notation, where  $n$  is the dimension of the problem and  $m$  is the number of examples.

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

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

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

**Smoothness and convexity.** A differentiable function  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  is called *convex* if for any  $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$  we have  $f(\mathbf{y}) \geq f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle$ . Furthermore,  $f$  is called  $\beta$ -*smooth* (with respect to some norm  $\|\cdot\|$ ) for some real number  $\beta > 0$  if for any  $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$  we have  $f(\mathbf{y}) \leq f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle + (\beta/2) \|\mathbf{y} - \mathbf{x}\|^2$ . If  $f$  is only  $\beta$ -smooth along  $s$ -sparse directions (i.e. only for  $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$  such that  $\|\mathbf{y} - \mathbf{x}\|_0 \leq s$ ), then we call  $f$   $\beta$ -smooth at *sparsity level*  $s$  and denote the smallest such  $\beta$  by  $\beta_s$  and call it the *restricted smoothness constant* (at sparsity level  $s$ ).

### 3 LOGISTIC REGRESSION ANALYSIS VIA MULTIPLICATIVE SMOOTHNESS

In the logistic regression problem, our goal is to minimize the function  $f(\mathbf{x}) = \sum_{i=1}^m \log(1 + e^{-(\mathbf{A}\mathbf{x})_i})$ ,

where  $\mathbf{A} \in \mathbb{R}^{m \times n}$  is a data matrix<sup>1</sup>

Our starting point, as is usually the case with first-order methods, will be the second order Taylor expansion of  $f$ :

$$f(\mathbf{x} + \tilde{\mathbf{x}}) = f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \tilde{\mathbf{x}} \rangle + \frac{1}{2} \langle \tilde{\mathbf{x}}, \nabla^2 f(\tilde{\mathbf{x}}) \tilde{\mathbf{x}} \rangle, \quad (1)$$

where, by the mean value theorem for twice continuously differentiable functions,  $\tilde{\mathbf{x}}$  is entry-wise between  $\mathbf{x}$  and  $\mathbf{x}'$ , and  $\nabla^2 f(\tilde{\mathbf{x}})$  is the Hessian of  $f$  at  $\tilde{\mathbf{x}}$ . In fact, as long as the step  $\tilde{\mathbf{x}}$  is not too large, the Hessian at  $\tilde{\mathbf{x}}$  will not differ much (spectrally) from the Hessian at  $\mathbf{x}$ . This is because of the following property of the logistic function called *second order robustness* Cohen et al. (2017), which is also very closely related to quasi-self-concordance Bach (2010).

**Definition 3.1** (Second-order robustness). *A twice differentiable function  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  is called  $q$ -second order robust with respect to a norm  $\|\cdot\|$  if its Hessian is stable in any  $(1/q)$ -sized  $\|\cdot\|$ -ball, i.e. for any  $\mathbf{x}, \mathbf{x}' \in \mathbb{R}^n$  such that  $\|\mathbf{x}' - \mathbf{x}\| \leq 1/q$ , we have  $\frac{1}{2} \nabla^2 f(\mathbf{x}) \preceq \nabla^2 f(\mathbf{x}') \preceq 2 \nabla^2 f(\mathbf{x})$ .*

It is not hard to see that  $f$  is  $2M$ -second order robust with respect to the  $\ell_1$  norm, where  $M$  is an upper bound on the entries of  $\mathbf{A}$  in absolute value. Because of this, (1) implies the much simpler

$$f(\mathbf{x} + \tilde{\mathbf{x}}) = f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \tilde{\mathbf{x}} \rangle + \langle \tilde{\mathbf{x}}, \nabla^2 f(\mathbf{x}) \tilde{\mathbf{x}} \rangle, \quad (2)$$

as long as  $\|\tilde{\mathbf{x}}\|_1 \leq 1/(2M)$ . We can easily calculate that  $\nabla f(\mathbf{x}) = -\mathbf{A}^\top (\mathbf{1} - \sigma(\mathbf{A}\mathbf{x}))$ , where  $\sigma(t) = 1/(1 + e^{-t})$  is the sigmoid function, and  $\nabla^2 f(\mathbf{x}) = \mathbf{A}^\top \text{diag}(\mathbf{w}(\mathbf{x})) \mathbf{A}$ , where  $\mathbf{w}(\mathbf{x}) =$

<sup>1</sup>This formulation is without loss of generality, because we can incorporate the binary  $\pm 1$  labels into the matrix  $\mathbf{A}$  and assume that all the labels are positive.

$\sigma(\mathbf{x})(1 - \sigma(\mathbf{x}))$  are diagonal weights. Now, we should note that the second order term of (1) can be re-written as  $\frac{1}{2}\langle \mathbf{w}(\mathbf{x}), (\mathbf{A}\tilde{\mathbf{x}})^2 \rangle$ . This term, whose magnitude is what will determine the step size of the algorithm and in turn the bound on the total number of iterations, becomes smaller as the weights  $\mathbf{w}(\mathbf{x})$  become smaller. The crucial observation is that these weights are bounded in a way that depends on the *loss* of  $\tilde{\mathbf{x}}$ , concretely:

$$\sum_{i=1}^m (\mathbf{w}(\mathbf{x}))_i \leq f(\mathbf{x}). \quad (3)$$

In other words, as the loss decreases,  $f$  becomes *smoother* (in an appropriate sense). This is the main observation on which our analysis is based, and is what allows the algorithm to employ a step size that is *inversely proportional* to the loss.

**Multiplicative smoothness.** The above discussion motivates the following definition of *multiplicative smoothness*. This is related to the usual definition of smoothness but also incorporates the property that the function becomes smoother as the loss decreases.

**Definition 3.2** (Multiplicative smoothness). *We call a twice differentiable function  $f : \mathbb{R}^n \rightarrow \mathbb{R}_{>0}$   $\mu$ -multiplicatively smooth with respect to a norm  $\|\cdot\|$ , if for any  $\mathbf{x}, \tilde{\mathbf{x}} \in \mathbb{R}^n$  we have*

$$\frac{\tilde{\mathbf{x}}^\top \nabla^2 f(\mathbf{x}) \tilde{\mathbf{x}}}{f(\mathbf{x})} \leq \mu \|\tilde{\mathbf{x}}\|^2.$$

Our use of a general norm is not an over-generalization, since as we will see the  $\ell_1$  norm is more suitable for sparse logistic regression, and the  $\ell_2$  norm is more suitable for the unrestricted case. In fact, it can be proved that  $f$  is  $M^2$ -multiplicatively smooth with respect to the  $\ell_1$  norm, where we remind that  $M$  is a bound on the entries of  $\mathbf{A}$  in absolute value.

In the following sections, we will see how the second order robustness and multiplicative smoothness properties play into the design and analysis of algorithms for sparse and general logistic regression.

## 4 SPARSE LOGISTIC REGRESSION

As we saw, the logistic loss is  $2M$ -second order robust and  $M^2$ -multiplicatively smooth with respect to the  $\ell_1$  norm. This is an ideal norm for *sparse* logistic regression, where in addition to minimizing the loss we want to restrict the solution to have few non-zero entries. In particular, it yields a variant of the  $\ell_1$  gradient descent algorithm (aka greedy coordinate descent), which is presented in Algorithm 1.

---

### Algorithm 1 Greedy Coordinate Descent

---

```

1: procedure GREEDYCOORDINATEDESCENT( $\mathbf{x}^0, T, M, B$ )
2:   Let  $f(\mathbf{x}) := \sum_{i=1}^m \log(1 + e^{-b_i(\mathbf{A}\mathbf{x})_i})$ 
3:   for  $t = 0, \dots, T - 1$  do
4:     For all  $i \in [n]$  define  $\zeta_i = \begin{cases} \lambda_t & \text{if } x_i^t = 0 \\ 0 & \text{if } \|\mathbf{x}^t\|_1 \geq B \text{ and } \nabla_i f(\mathbf{x}^t) \cdot x_i^t < 0 \\ 1 & \text{otherwise} \end{cases}$ 
5:      $i \leftarrow \operatorname{argmax}_i \{\zeta_i |\nabla_i f(\mathbf{x}^t)|\}$ 
6:      $\eta \leftarrow (2M \max\{Mf(\mathbf{x}^t), |\nabla_i f(\mathbf{x}^t)|\})^{-1}$ 
7:      $x_i^{t+1} \leftarrow x_i^t - \eta \nabla_i f(\mathbf{x}^t)$ 
return  $\mathbf{x}^T$ 

```

---

The first thing that should be noted about this algorithm is the crucial parameters  $\lambda_t$ . These parameters offer a quantitative threshold between sparsity and speed of convergence. In particular, when  $\lambda_t$  is 1, then all entries (regardless of whether they are zero or not) are treated the same. When  $\lambda_t \ll 1$ , on the other hand, the gradient entries corresponding to zero entries are discounted by a factor  $\ll 1$ , thus making the algorithm less eager to update these as opposed to non-zero entries, whose update doesn't increase sparsity.

A practical consideration about Algorithm 1 is order. The second condition in line 4 is to make sure that the  $\ell_1$  norm of the solution never exceeds a given bound on the  $\ell_1$  norm of the optimal solution. This check is useful for the theoretical analysis but should likely be removed in any practical implementation.

We are ready for the main theorem of this section. In the proof, which can be found in Appendix A.2.2, we present an analysis of Algorithm 1 for sparse logistic regression. In addition to an upper bound  $B \geq \|\mathbf{x}^*\|_\infty$ , it also requires an approximation  $B_1$  of  $\|\mathbf{x}^*\|_1$ . One possible approach is to approximate it by  $B$ , but in practice this would be a learning rate hyperparameter to be tuned.

**Theorem 4.1** (Sparse logistic regression). *Given a binary classification instance ( $\mathbf{A} \in [-M, M]^{m \times n}$ ,  $\mathbf{b} \in \{1, -1\}^m$ ) and for any solution  $\mathbf{x}^* \in [-B, B]^n$  with  $M \geq \max\{\|\mathbf{x}^*\|_\infty^{-1}, B^{-1}\}^2$  and a known parameter  $B_1 \in [\frac{1}{C}\|\mathbf{x}^*\|_1, \|\mathbf{x}^*\|_1]$  for some  $C \geq 1$ , Algorithm 1 with  $\lambda_t = \min\{B_1/\|\mathbf{x}^t\|_1, 1\}$ , initial solution  $\mathbf{x}^0 \in \mathbb{R}^n$ , and error tolerance  $0 < \varepsilon < m/2$  returns a solution  $\mathbf{x}$  with*

$$f(\mathbf{x}) \leq (1 + \delta) \cdot f(\mathbf{x}^*) + \varepsilon$$

and sparsity

$$s' := \|\mathbf{x}\|_0 = O\left(\|\mathbf{x}^*\|_1^2 M^2 \left(\frac{1}{\delta} + \log \frac{f(\mathbf{x}^0) - f(\mathbf{x}^*)}{\varepsilon}\right)\right)$$

in

$$T = O\left(\left(\|\mathbf{x}\|_0^2 + \|\mathbf{x}^*\|_0^2\right) M^2 B^2 C^2 \left(\frac{1}{\delta} + \log \frac{f(\mathbf{x}^0) - f(\mathbf{x}^*)}{\varepsilon}\right)\right)$$

iterations, for any choice of  $\delta \in (0, 1)$  and parameter  $c > 0$ . Each iteration consists of evaluating the logistic regression gradient  $\nabla f$  plus  $O(m + n)$  additional time.

**Corollary 4.2.** *If  $M, B, C \leq \tilde{O}(1)$  and  $\mathbf{x}^*$  is  $s$ -sparse, then Algorithm 1 with  $\lambda_t = \min\{1/\|\mathbf{x}^t\|_1, 1\}$  returns a solution  $\mathbf{x}$  with*

$$f(\mathbf{x}) \leq 1.1 \cdot f(\mathbf{x}^*) + \varepsilon$$

and sparsity

$$s' := \|\mathbf{x}\|_0 = \tilde{O}\left(s^2 \log \frac{1}{\varepsilon}\right)$$

in

$$T = \tilde{O}\left(s^4 \log^3 \frac{1}{\varepsilon}\right)$$

iterations.

It is useful to compare these results to the results of Shalev-Shwartz et al. (2010) for sparse optimization of general smooth convex functions. Even though they achieve the stronger error bound of  $f(\mathbf{x}) \leq f(\mathbf{x}^*) + \varepsilon$ , the sparsity of the final solution is in the order of  $s^2 \frac{m}{\varepsilon}$ , which has an exponentially worse error dependence than  $s^2 \log \frac{m}{\varepsilon}$ . Therefore, if the approximation rate  $(1 + \delta)$  is tolerable in front of  $f(\mathbf{x}^*)$ , then one can obtain exponentially faster sparsity and convergence.

If we are willing to perform fully corrective steps as described in Algorithm 2, then we can get a cleaner and slightly simpler analysis. This is presented in Theorem 4.3 and proved in Appendix A.2.3. Fully corrective steps can be useful when there is an efficient (dense) optimization algorithm and one wishes to use it as a black box for sparse optimization. In practice, one does not need to perform a full correction, but only a small number of corrective (usually gradient) steps over the current support of the solution.

---

<sup>2</sup>the theorem can be stated without this additional constraint, but we include it because it makes the bounds considerably simpler

**Algorithm 2** Greedy coordinate descent with fully corrective steps

---

```

1: procedure FULLYCORRECTIVEGREEDYCOORDINATEDESCENT( $\mathbf{x}^0, T, M, B$ )
2:   Let  $f(\mathbf{x}) := \sum_{i=1}^m \log(1 + e^{-b_i(\mathbf{A}\mathbf{x})_i})$ 
3:    $S^0 \leftarrow \text{supp}(\mathbf{x}^0)$ 
4:   for  $t = 0, \dots, T - 1$  do
5:      $i \leftarrow \text{argmax}_i \{|\nabla_i f(\mathbf{x}^t)|\}$ 
6:      $S^{t+1} \leftarrow S^t \cup \{i\}$ 
7:      $\mathbf{x}^{t+1} \leftarrow \underset{\mathbf{x}: \text{supp}(\mathbf{x}) \subseteq S^{t+1}}{\text{argmin}} f(\mathbf{x})$ 
   return  $\mathbf{x}^T$ 

```

---

**Theorem 4.3** (Sparse logistic regression with fully corrective steps). *Given a binary classification instance ( $\mathbf{A} \in [-M, M]^{m \times n}$ ,  $\mathbf{b} \in \{1, -1\}^m$ ) and for any solution  $\mathbf{x}^* \in \mathbb{R}^n$ , Algorithm 2 with error tolerance  $0 < \varepsilon < m/2$  and initial solution  $\mathbf{x}^0$  returns a solution  $\mathbf{x}$  with*

$$f(\mathbf{x}) \leq (1 + \delta) \cdot f(\mathbf{x}^*) + \varepsilon$$

and sparsity

$$s' := \|\mathbf{x}\|_0 = \|\mathbf{x}^0\|_0 + O\left(\|\mathbf{x}^*\|_1^2 M^2 \left(\frac{1}{\delta} + \log \frac{f(\mathbf{x}^0) - f(\mathbf{x}^*)}{\varepsilon}\right)\right)$$

in  $T = \|\mathbf{x}\|_0$  iterations, for any choice of  $\delta \in (0, 1)$ . Each iteration consists of evaluating the logistic regression gradient  $\nabla f$ , solving a logistic regression problem on  $s'$  variables, plus  $O(m+n)$  additional time.

## 5 DENSE LOGISTIC REGRESSION

In this section, our goal is to minimize the logistic function  $f$  without any constraint on the sparsity of the solution. The results of the Section 4 applied to a full sparsity of  $n$  already imply Corollary 5.1.

**Corollary 5.1** (Dense logistic regression). *Given a binary classification instance ( $\mathbf{A} \in [-M, M]^{m \times n}$ ,  $\mathbf{b} \in \{-1, 1\}^m$ ) and for any solution  $\mathbf{x}^* \in [-B, B]^n$  with  $M \geq \max\{\|\mathbf{x}^*\|_\infty^{-1}, B^{-1}\}$ , Algorithm 1 with  $\lambda_t = 1$  for all  $t$ , initial solution  $\mathbf{x}^0 \in \mathbb{R}^n$ , and error tolerance  $0 < \varepsilon < m/2$  returns a solution  $\mathbf{x}$  with*

$$f(\mathbf{x}) \leq (1 + \delta) \cdot f(\mathbf{x}^*) + \varepsilon$$

in

$$T = O\left(n^2 M^2 B^2 \left(\frac{1}{\delta} + \log \frac{f(\mathbf{x}^0) - f(\mathbf{x}^*)}{\varepsilon}\right)\right).$$

iterations, for any choice of  $\delta \in (0, 1)$ . Additionally,  $\|\mathbf{x}\|_\infty \leq B + \frac{1}{2M}$ . Each iteration consists of evaluating the logistic regression gradient  $\nabla f$  plus  $O(m+n)$  additional time.

Even though Corollary 5.1 has the same favorable convergence in terms of  $\delta$  and  $\varepsilon$  as Theorem 4.1, based on practical intuitions we would expect ( $\ell_2$ -based) gradient descent to perform better than greedy coordinate descent, which only updates one coordinate at a time, while having access to the full gradient. In fact, we can verify that the logistic loss does have the multiplicative smoothness condition with respect to the  $\ell_2$  norm, albeit in an almost trivial sense:

$$\langle \mathbf{w}(\mathbf{x}), (\mathbf{A}\mathbf{x})^2 \rangle \leq \|\mathbf{w}(\mathbf{x})\|_1 \|\mathbf{A}\mathbf{x}\|_\infty^2 \leq f(\mathbf{x}) \|\mathbf{A}\|_{2 \rightarrow \infty}^2 \|\mathbf{x}\|_2^2 \leq f(\mathbf{x}) \beta \|\mathbf{x}\|_2^2.$$

Here, using the inequality  $\|\mathbf{A}\|_{2 \rightarrow \infty}^2 \leq \|\mathbf{A}\|_2^2 := \beta$  implies  $\beta$ -multiplicative smoothness with respect to the  $\ell_2$  norm. Unfortunately, this is not significantly better than the  $\ell_1$  case: The number of iterations will be proportional to  $\beta \|\mathbf{x}^*\|_2^2$ , which can be  $\gg m$ .

Table 3: Upper bounds on the quantity  $\langle \mathbf{w}(\mathbf{x}), (\mathbf{A}\nabla f(\mathbf{x}))^2 \rangle / (f(\mathbf{x})m^{-1} \|\mathbf{A}\nabla f(\mathbf{x})\|_2^2)$ . Shown here is the maximum of this over  $\mathbf{x}$  being one of the first 1000 iterates.

Dataset	Max ratio	Dataset	Max ratio	Dataset	Max ratio
letter	0.40	skin	0.44	census	0.50
rcv1.test	0.36	w8all	0.40	adult	0.40
ijcnn1	0.47	shuttle.binary	0.37	poker	0.36
vehv2binary	0.37	kddcup04.phy	0.36	nomao	0.50
magic04	0.37	kddcup04.bio	0.48	covtype	0.36

Interestingly, real logistic regression instances exhibit the  $\ell_2$  multiplicative smoothness property with significantly better constants. In our experiments we found that along the path of gradients encountered by gradient descent in a variety of instances, the following property was true:

$$\langle \mathbf{w}(\mathbf{x}), (\mathbf{A}\nabla f(\mathbf{x}))^2 \rangle \leq f(\mathbf{x})\beta m^{-1} \|\nabla f(\mathbf{x})\|_2^2$$

This is an *effective*  $\beta m^{-1}$ -multiplicative smoothness property, because it is only assumed to be true for  $\mathbf{x}$ 's encountered by the gradient descent algorithm. As such, it is an empirical property. In order to check our hypothesis, we have run the gradient descent algorithm with the step sizes that are implied by Theorem 5.2, which we will see later. For each of the 15 experiments, we have run gradient descent for 1000 iterations, and calculated the maximum of the following quantity, over all iterations:

$$\frac{\langle \mathbf{w}(\mathbf{x}), (\mathbf{A}\nabla f(\mathbf{x}))^2 \rangle}{f(\mathbf{x})m^{-1} \|\mathbf{A}\nabla f(\mathbf{x})\|_2^2}.$$

If this is bounded by 1, and using the fact that  $\|\mathbf{A}\nabla f(\mathbf{x})\|_2^2 \leq \beta \|\nabla f(\mathbf{x})\|_2^2$ , this implies that  $f$  is effectively  $\beta m^{-1}$ -multiplicatively smooth with respect to the  $\ell_2$  norm. Indeed, as we can see in Table 3, these values are indeed less than 1 for all datasets and all iterations.

In the following, our plan is to prove convergence, *assuming* that  $f$  has the multiplicative smoothness property with the constants in our hypothesis above. Under this assumption, we can now prove a much stronger convergence theorem (here we are also using the fact that  $M^2 \leq \beta$  to replace  $2M$ -by  $2\sqrt{\beta}$ -second order robustness):

**Theorem 5.2.** *Let  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  be a convex function that is  $2\sqrt{\beta}$ -second order robust with respect to the  $\ell_1$  norm and  $\beta m^{-1}$ -multiplicatively smooth with respect to the  $\ell_2$  norm. Let  $\mathbf{x}^0 \in \mathbb{R}^n$  be an initial solution and  $\mathbf{x}^* \in \mathbb{R}^n$  be an arbitrary solution, where  $R := \|\mathbf{x}^0 - \mathbf{x}^*\|_2$  and  $R \geq \sqrt{n}$ .<sup>3</sup> Then, gradient descent with step size  $\eta_t = 0.5 \min \left\{ \frac{1}{\beta m^{-1} f(\mathbf{x})}, \frac{1}{\sqrt{\beta} \|\nabla f(\mathbf{x})\|_1} \right\}$  returns a solution with*

$$f(\mathbf{x}) \leq (1 + \delta)f(\mathbf{x}^*) + \varepsilon$$

after

$$T = O \left( \frac{\beta R^2}{m} \left( \frac{1}{\delta} + \log \frac{f(\mathbf{x}^0) - f(\mathbf{x}^*)}{\varepsilon} \right) \right)$$

iterations.

## 6 MAXIMUM MARGIN SOLUTIONS

It is known that running gradient descent on the logistic loss on linearly separable data converges to the hard SVM (maximum margin) classifier Soudry et al. (2018), yet at the slow rate of

$$\left\| \frac{\mathbf{x}^t}{\|\mathbf{x}^t\|_2} - \frac{\mathbf{x}^*}{\|\mathbf{x}^*\|_2} \right\|_2 \leq O \left( \frac{\log \log t}{\log t} \right).$$

<sup>3</sup>The assumption  $R \geq \sqrt{n}$  is not necessary but simplifies the bounds.



As all our results work best when the data is separable, it is natural to ask about what they imply for margin maximization.

We consider the constrained logistic regression problem

$$\min_{\|\mathbf{x}\|_2 \leq 1} f_p(\mathbf{x}) := \sum_i \log(1 + e^{-pb_i(\mathbf{A}\mathbf{x})_i}) \quad (4)$$

We start by observing that Corollary 5.1 and Theorem 5.2 can be modified to solve (4), with a blowup of  $p^2$  in the number of iterations. In particular, the number of iterations will be  $\tilde{O}\left(p^2 X \left(\frac{1}{\delta} + \log \frac{1}{\varepsilon}\right)\right)$ , where  $X$  depends on whether we use Corollary 5.1 or Theorem 5.2, but is beyond the point of this section, since here we are interested in the error dependence.

Picking  $\delta = 1$ ,  $\varepsilon = me^{-p\alpha}$ , and  $p = \frac{\log(6m)}{\alpha\hat{\varepsilon}}$  for some target error  $\hat{\varepsilon} \in (0, 1)$ , we get the following theorem:

**Theorem 6.1.** *Consider a linearly separable binary classification instance ( $\mathbf{A} \in \mathbb{R}^{m \times n}$ ,  $\mathbf{b} \in \{1, -1\}^m$ ), and a solution  $\mathbf{x}^*$  that maximizes  $\min_i \frac{b_i(\mathbf{A}\mathbf{x}^*)_i}{\|\mathbf{x}^*\|_2} := \alpha$ . Then, we can obtain a solution  $\mathbf{x}$  with  $\|\mathbf{x}\|_2 \leq 1$  and  $f_p(\mathbf{x}) \leq 3f_p(\mathbf{x}^*)$  in  $\tilde{O}\left(X \frac{1}{\alpha^2 \varepsilon^3}\right)$  iterations of gradient descent, where  $p = \frac{\log(6m)}{\alpha\hat{\varepsilon}}$ . Furthermore,  $\mathbf{x}$  has  $(1 - \hat{\varepsilon})$ -optimal margins:*

$$\min_i \frac{b_i(\mathbf{A}\mathbf{x})_i}{\|\mathbf{x}\|_2} \geq \alpha(1 - \hat{\varepsilon})$$

and is close to the maximum margin classifier:

$$\left\| \frac{\mathbf{x}}{\|\mathbf{x}\|_2} - \frac{\mathbf{x}^*}{\|\mathbf{x}^*\|_2} \right\|_2 \leq 2\sqrt{\hat{\varepsilon}}$$

It is not hard to see that Theorem 6.1 gives an exponential improvement in the error dependence compared to Soudry et al. (2018).

## 7 NUMERICAL EXAMPLE

In order to numerically validate our algorithm, we run logistic regression on the well known UCI adult binary classification dataset. In order to simulate a separable dataset, we first run gradient descent on the whole data, and then discard the misclassified data points. This gives us a separable dataset. Then, we run two variants of gradient descent: One with constant step size given by  $\beta^{-1}$ , and one with increasing step size given by  $\eta_t = \beta^{-1} f(\mathbf{x}^0)/f(\mathbf{x}^t)$ , with no other change. This is motivated by our findings, which suggest that the step size should increase proportionally to the decrease of the loss. As we can see in Figure 2, the error in the case of fixed step size decays as  $\text{poly}(1/t)$ , while in the case of increasing step size we have linear convergence (albeit with a low rate because the margins are in the order of  $10^{-6}$ ).

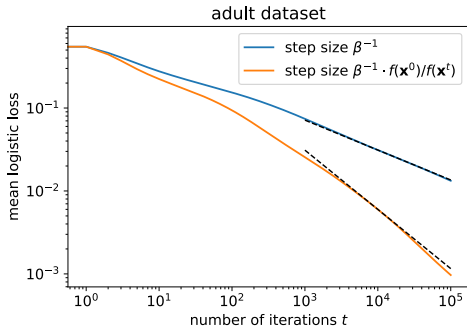


Figure 2: Comparison of fixed vs increasing step size on logistic regression on adult dataset

## REFERENCES

- Deeksha Adil, Brian Bullins, and Sushant Sachdeva. Unifying width-reduced methods for quasi-self-concordant optimization. *Advances in Neural Information Processing Systems*, 34, 2021.
- Kyriakos Axiotis and Maxim Sviridenko. Sparse convex optimization via adaptively regularized hard thresholding. *Journal of Machine Learning Research*, 22:1–47, 2021.
- Kyriakos Axiotis and Maxim Sviridenko. Iterative hard thresholding with adaptive regularization: Sparser solutions without sacrificing runtime. *arXiv preprint arXiv:2204.08274*, 2022.
- Francis Bach. Self-concordant analysis for logistic regression. *Electronic Journal of Statistics*, 4: 384–414, 2010.
- Thomas Blumensath and Mike E Davies. Iterative hard thresholding for compressed sensing. *Applied and computational harmonic analysis*, 27(3):265–274, 2009.
- Florentina Bunea. Honest variable selection in linear and logistic regression models via  $\ell_1$  and  $\ell_1 + \ell_2$  penalization. *Electronic Journal of Statistics*, 2:1153–1194, 2008.
- Yair Carmon, Arun Jambulapati, Qijia Jiang, Yujia Jin, Yin Tat Lee, Aaron Sidford, and Kevin Tian. Acceleration with a ball optimization oracle. *Advances in Neural Information Processing Systems*, 33:19052–19063, 2020.
- Michael B Cohen, Aleksander Madry, Dimitris Tsipras, and Adrian Vladu. Matrix scaling and balancing via box constrained newton’s method and interior point methods. In *2017 IEEE 58th Annual Symposium on Foundations of Computer Science (FOCS)*, pp. 902–913. IEEE, 2017.
- Robert M Freund, Paul Grigas, and Rahul Mazumder. Condition number analysis of logistic regression, and its implications for standard first-order solution methods. *arXiv preprint arXiv:1810.08727*, 2018.
- Prateek Jain, Ambuj Tewari, and Inderjit S Dhillon. Orthogonal matching pursuit with replacement. In *Advances in neural information processing systems*, pp. 1215–1223, 2011.
- Prateek Jain, Ambuj Tewari, and Purushottam Kar. On iterative hard thresholding methods for high-dimensional m-estimation. In *Advances in Neural Information Processing Systems*, pp. 685–693, 2014.
- Ziwei Ji and Matus Telgarsky. Risk and parameter convergence of logistic regression. *arXiv preprint arXiv:1803.07300*, 2018.
- Sai Praneeth Karimireddy, Sebastian U Stich, and Martin Jaggi. Global linear convergence of newton’s method without strong-convexity or lipschitz gradients. *arXiv preprint arXiv:1806.00413*, 2018.
- Mor Shpigel Nacson, Jason Lee, Suriya Gunasekar, Pedro Henrique Pamplona Savarese, Nathan Srebro, and Daniel Soudry. Convergence of gradient descent on separable data. In *The 22nd International Conference on Artificial Intelligence and Statistics*, pp. 3420–3428. PMLR, 2019.
- Balas Kausik Natarajan. Sparse approximate solutions to linear systems. *SIAM journal on computing*, 24(2):227–234, 1995.
- Gunnar Rätsch, Sebastian Mika, and Manfred KK Warmuth. On the convergence of leveraging. *Advances in Neural Information Processing Systems*, 14, 2001.
- Shai Shalev-Shwartz, Nathan Srebro, and Tong Zhang. Trading accuracy for sparsity in optimization problems with sparsity constraints. *SIAM Journal on Optimization*, 20(6):2807–2832, 2010.
- Daniel Soudry, Elad Hoffer, Mor Shpigel Nacson, Suriya Gunasekar, and Nathan Srebro. The implicit bias of gradient descent on separable data. *The Journal of Machine Learning Research*, 19(1):2822–2878, 2018.
- Matus Telgarsky. Margins, shrinkage, and boosting. In *International Conference on Machine Learning*, pp. 307–315. PMLR, 2013.

Matus Telgarsky and Yoram Singer. A primal-dual convergence analysis of boosting. *Journal of Machine Learning Research*, 13(3), 2012.

Sara A Van de Geer. High-dimensional generalized linear models and the lasso. *The Annals of Statistics*, 36(2):614–645, 2008.

## A MISSING PROOFS FROM SECTION 4

### A.1 PROOF OF MAIN LEMMA

**Lemma A.1** (Gradient lower bound). *Let  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  be a differentiable convex function and let  $\mathbf{x} \in [-B', B']^n$ ,  $\mathbf{x}^* \in [-B, B]^n$  be two solutions for some parameters  $B' \geq B > 0$ . For all  $i \in [n]$  we define*

$$\zeta_i = \begin{cases} \lambda & \text{if } x_i = 0 \\ 0 & \text{if } |x_i| \geq B \text{ and } \nabla_i f(\mathbf{x}) \cdot x_i < 0 \\ 1 & \text{otherwise} \end{cases}$$

where  $0 < \lambda \leq 1$ , and let  $i^* = \operatorname{argmax}_i \{\zeta_i |\nabla_i f(\mathbf{x})|\}$ . Then, at least one of the following is true:

- $|\nabla_{i^*} f(\mathbf{x})| \geq \frac{f(\mathbf{x}) - f(\mathbf{x}^*)}{\|\mathbf{x}^*\|_1 + \lambda \|\mathbf{x}\|_1}$
- $|\nabla_{i^*} f(\mathbf{x})| \geq \frac{f(\mathbf{x}) - f(\mathbf{x}^*)}{\lambda^{-1} \|\mathbf{x}^*\|_1 + \|\mathbf{x}\|_1}$  and  $x_i \neq 0$ .

*Proof.* Let  $S = \{i \mid x_i \neq 0\}$  and  $F = \{i \mid |x_i| < B \text{ or } \nabla_i f(\mathbf{x}) \cdot x_i \geq 0\}$ . By convexity of  $f$ , we have

$$\begin{aligned} f(\mathbf{x}^*) &\geq f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \mathbf{x}^* - \mathbf{x} \rangle \\ &\geq f(\mathbf{x}) + \langle \nabla_F f(\mathbf{x}), \mathbf{x}^* - \mathbf{x} \rangle \\ &= f(\mathbf{x}) + \langle \nabla_F f(\mathbf{x}), \mathbf{x}^* \rangle - \langle \nabla_{S \cap F} f(\mathbf{x}), \mathbf{x} \rangle \\ &\geq f(\mathbf{x}) - \|\nabla_F f(\mathbf{x})\|_\infty \|\mathbf{x}^*\|_1 - \|\nabla_{S \cap F} f(\mathbf{x})\|_\infty \|\mathbf{x}\|_1, \end{aligned}$$

therefore

$$\|\nabla_F f(\mathbf{x})\|_\infty \|\mathbf{x}^*\|_1 + \|\nabla_{S \cap F} f(\mathbf{x})\|_\infty \|\mathbf{x}\|_1 \geq f(\mathbf{x}) - f(\mathbf{x}^*). \quad (5)$$

Now, if  $i^* \notin S$ , by definition of the  $\zeta_i$ 's and  $i^*$  we have

$$\lambda \|\nabla_F f(\mathbf{x})\|_\infty = \lambda |\nabla_{i^*} f(\mathbf{x})| \geq \|\nabla_{S \cap F} f(\mathbf{x})\|_\infty.$$

and so (5) implies

$$\begin{aligned} |\nabla_{i^*} f(\mathbf{x})| \|\mathbf{x}^*\|_1 + \lambda |\nabla_{i^*} f(\mathbf{x})| \|\mathbf{x}\|_1 &\geq f(\mathbf{x}) - f(\mathbf{x}^*) \\ \Rightarrow |\nabla_{i^*} f(\mathbf{x})| &\geq \frac{f(\mathbf{x}) - f(\mathbf{x}^*)}{\|\mathbf{x}^*\|_1 + \lambda \|\mathbf{x}\|_1}. \end{aligned}$$

Otherwise if  $i^* \in S$ , we have

$$\|\nabla_{S \cap F} f(\mathbf{x})\|_\infty = |\nabla_{i^*} f(\mathbf{x})| \geq \lambda \|\nabla_F f(\mathbf{x})\|_\infty,$$

and so (5) implies

$$\begin{aligned} \lambda^{-1} |\nabla_{i^*} f(\mathbf{x})| \|\mathbf{x}^*\|_1 + |\nabla_{i^*} f(\mathbf{x})| \|\mathbf{x}\|_1 &\geq f(\mathbf{x}) - f(\mathbf{x}^*) \\ \Rightarrow |\nabla_{i^*} f(\mathbf{x})| &\geq \frac{f(\mathbf{x}) - f(\mathbf{x}^*)}{\lambda^{-1} \|\mathbf{x}^*\|_1 + \|\mathbf{x}\|_1}. \end{aligned}$$

□

**Lemma A.2** (Coordinate update). *Let  $f : \mathbb{R}^n \rightarrow \mathbb{R}_{\geq 0}$  be a twice continuously differentiable convex function that is  $2\gamma$ -second order robust and  $\gamma^2$ -multiplicatively smooth with respect to the  $\ell_1$  norm, for some  $\gamma > 0$ . Let  $\mathbf{x} \in [-B', B']^n$  be a suboptimal solution such that  $f(\mathbf{x}) \geq f(\mathbf{x}^*)$ , where  $\mathbf{x}^* \in [-B, B]^n$  is some unknown solution with  $\gamma \|\mathbf{x}^*\|_1 \geq 1$ , and  $B' \geq B > 0$  are some parameters. We make the update*

$$\mathbf{x}' = \mathbf{x} - \eta \nabla_i f(\mathbf{x}) \mathbf{1}_i,$$

where  $i$  is picked as in Lemma A.1 for some parameter  $\lambda \in (0, 1)$  and  $\eta = 0.5 \min \left\{ \frac{1}{\gamma^2 f(\mathbf{x})}, \frac{1}{\gamma |\nabla_i f(\mathbf{x})|} \right\}$  is a step size. Then, at least one of the following is true about the progress in decreasing  $f$ :

- $f(\mathbf{x}) - f(\mathbf{x}') \geq \frac{(f(\mathbf{x}) - f(\mathbf{x}^*))^2}{4\gamma^2 f(\mathbf{x}) (\|\mathbf{x}^*\|_1 + \lambda \|\mathbf{x}\|_1)^2}$
- $f(\mathbf{x}) - f(\mathbf{x}') \geq \frac{(f(\mathbf{x}) - f(\mathbf{x}^*))^2}{4\gamma^2 f(\mathbf{x}) (\lambda^{-1} \|\mathbf{x}^*\|_1 + \|\mathbf{x}\|_1)^2}$  and  $x_i \neq 0$ ,

and the norm of the new solution is bounded as  $\|\mathbf{x}'\|_\infty \leq \max\{B', B + \frac{1}{2\gamma}\}$ . In the case that  $f(\mathbf{x}) < f(\mathbf{x}^*)$  we have  $f(\mathbf{x}') \leq f(\mathbf{x})$ .

*Proof.* We first consider a generic update  $\mathbf{x}' = \mathbf{x} + \tilde{\mathbf{x}}$ . By Taylor's theorem and the fact that  $f$  is twice continuously differentiable, we have

$$f(\mathbf{x}') = f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \tilde{\mathbf{x}} \rangle + \frac{1}{2} \langle \tilde{\mathbf{x}}, \nabla^2 f(\bar{\mathbf{x}}) \tilde{\mathbf{x}} \rangle,$$

for some  $\bar{\mathbf{x}}$  that is entrywise between  $\mathbf{x}$  and  $\tilde{\mathbf{x}}$ .

Since  $f$  is  $2\gamma$ -second-order-robust and  $\gamma^2$ -multiplicatively-smooth with respect to the  $\ell_1$  norm, as long as the update is bounded in  $\ell_1$  norm as

$$\|\tilde{\mathbf{x}}\|_1 \leq 1/(2\gamma) \tag{6}$$

we have

$$\begin{aligned} f(\mathbf{x}') &\leq f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \tilde{\mathbf{x}} \rangle + \langle \tilde{\mathbf{x}}, \nabla^2 f(\mathbf{x}) \tilde{\mathbf{x}} \rangle \\ &\leq f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \tilde{\mathbf{x}} \rangle + \gamma^2 f(\mathbf{x}) \|\tilde{\mathbf{x}}\|_1^2. \end{aligned}$$

Note that the right hand side is minimized for

$$\tilde{\mathbf{x}} = -\frac{H_1(\nabla f(\mathbf{x}))}{2\gamma^2 f(\mathbf{x})},$$

where  $H_1$  is the hard thresholding operator that zeroes out all but the top entry in absolute value. This is a coordinate descent step. Our step will be slightly more careful so that it doesn't unnecessarily increase the sparsity of  $\mathbf{x}$ . We consider the following coordinate step

$$\tilde{\mathbf{x}} = -\eta \nabla_i f(\mathbf{x}),$$

where  $\eta > 0$  and  $i$  are as defined in the lemma statement. We now have

$$f(\mathbf{x}) - f(\mathbf{x}') \geq (\eta - \eta^2 \gamma^2 f(\mathbf{x})) (\nabla_i f(\mathbf{x}))^2$$

The term  $(\eta - \eta^2 \gamma^2 f(\mathbf{x}))$  is maximized at  $\eta = \frac{1}{2\gamma^2 f(\mathbf{x})}$ . In addition, to stay in the  $\ell_1$  neighborhood where the Hessian is stable, we need to satisfy (6) by making sure that  $\eta \leq \frac{1}{2\gamma |\nabla_i f(\mathbf{x})|}$ . Based on these requirements, we pick

$$\eta = \min \left\{ \frac{1}{2\gamma^2 f(\mathbf{x})}, \frac{1}{2\gamma |\nabla_i f(\mathbf{x})|} \right\}$$

and conclude that

$$\begin{aligned} f(\mathbf{x}) - f(\mathbf{x}') &\geq \min \left\{ \frac{1}{4\gamma^2 f(\mathbf{x})}, \frac{1}{4\gamma |\nabla_i f(\mathbf{x})|} \right\} (\nabla_i f(\mathbf{x}))^2 \\ &= \min \left\{ \frac{(\nabla_i f(\mathbf{x}))^2}{4\gamma^2 f(\mathbf{x})}, \frac{|\nabla_i f(\mathbf{x})|}{4\gamma} \right\}. \end{aligned}$$

Note that this is always  $\geq 0$  and so we have  $f(\mathbf{x}') \leq f(\mathbf{x})$  even if  $f(\mathbf{x}) < f(\mathbf{x}^*)$ . We now take two cases and use the two bullets of Lemma A.1 accordingly.

**Case 1:**  $x_i = 0$ . The first bullet of Lemma A.1 has to be true, i.e.

$$|\nabla_i f(\mathbf{x})| \geq \frac{f(\mathbf{x}) - f(\mathbf{x}^*)}{\|\mathbf{x}^*\|_1 + \lambda \|\mathbf{x}\|_1}.$$

Therefore,

$$\begin{aligned} f(\mathbf{x}) - f(\mathbf{x}') &\geq \min \left\{ \frac{(f(\mathbf{x}) - f(\mathbf{x}^*))^2}{4\gamma^2 f(\mathbf{x}) (\|\mathbf{x}^*\|_1 + \lambda \|\mathbf{x}\|_1)^2}, \frac{f(\mathbf{x}) - f(\mathbf{x}^*)}{4\gamma (\|\mathbf{x}^*\|_1 + \lambda \|\mathbf{x}\|_1)} \right\} \\ &= \frac{(f(\mathbf{x}) - f(\mathbf{x}^*))^2}{4\gamma^2 f(\mathbf{x}) (\|\mathbf{x}^*\|_1 + \lambda \|\mathbf{x}\|_1)^2}, \end{aligned}$$

where we used the facts that  $f(\mathbf{x}) - f(\mathbf{x}^*) \leq f(\mathbf{x})$  and  $\gamma \|\mathbf{x}^*\|_1 \geq 1$ .

**Case 2:**  $x_i \neq 0$ . If the first bullet of Lemma A.1 is true, we can proceed as in the previous case. Otherwise, we use the second bullet of Lemma A.1 and similarly get

$$f(\mathbf{x}) - f(\mathbf{x}') \geq \frac{(f(\mathbf{x}) - f(\mathbf{x}^*))^2}{4\gamma^2 f(\mathbf{x}) (\lambda^{-1} \|\mathbf{x}^*\|_1 + \|\mathbf{x}\|_1)^2}.$$

Finally, in order to bound  $\|\mathbf{x}'\|_\infty$ , we first note that  $\|\mathbf{x}\|_\infty \leq B'$ . Now, by our choice of  $i$  we have that either  $|x_i| < B$ , or  $\nabla_i f(\mathbf{x}) \cdot x_i > 0$ . In the first case, we have

$$|x'_i| \leq |x_i| + |\tilde{x}_i| < B + \frac{1}{2\gamma},$$

where we used (6). Otherwise, we have that  $|x_i| \geq B$  and  $\nabla_i f(\mathbf{x}) \cdot x_i > 0$ . This implies that  $x_i$  and  $\tilde{x}_i$  have different signs, so

$$|x'_i| = |x_i + \tilde{x}_i| \leq \max\{|x_i|, |\tilde{x}_i|\} \leq \max\left\{B', \frac{1}{2\gamma}\right\}.$$

Therefore, in any case we have  $|x'_i| \leq \max\left\{B', B + \frac{1}{2\gamma}\right\}$ .  $\square$

## A.2 PROOFS OF THEOREMS

### A.2.1 PROOF OF COROLLARY 5.1

*Proof.* We will apply Lemma A.2 for  $T$  iterations to obtain solutions  $\mathbf{x}^0, \dots, \mathbf{x}^T$ , where for some  $T$  that will be defined later. The logistic function  $f$  is  $2M$ -second order robust and  $M^2$ -multiplicatively smooth with respect to the  $\ell_1$  norm, so Lemma A.2 can be applied with  $\gamma = M$  and  $B' = B + \frac{1}{2M}$ .

Based on the guarantee of Lemma A.2, we get the following bound on the  $\ell_1$  norm of  $\mathbf{x}^t$  at all times:

$$\|\mathbf{x}^t\|_1 \leq n \|\mathbf{x}^t\|_\infty \leq n \left( B + \frac{1}{2M} \right) \leq (3/2)nB.$$

Let  $\bar{t}$  be the smallest  $t \geq 0$  for which  $f(\mathbf{x}^{\bar{t}}) \leq 2f(\mathbf{x}^*)$  or  $f(\mathbf{x}^{\bar{t}}) \leq f(\mathbf{x}^*) + \varepsilon$ , and let  $\bar{t} = \infty$  if this never happens. Therefore, for all  $t < \bar{t}$  we have  $f(\mathbf{x}^t) \geq 2f(\mathbf{x}^*) \Rightarrow \frac{f(\mathbf{x}^t) - f(\mathbf{x}^*)}{f(\mathbf{x}^t)} \geq \frac{1}{2}$ , and so the statement of Lemma A.2 gives:

$$\begin{aligned} f(\mathbf{x}^t) - f(\mathbf{x}^{t+1}) &\geq \frac{f(\mathbf{x}^t) - f(\mathbf{x}^*)}{8M^2(\|\mathbf{x}^*\|_1 + \|\mathbf{x}^t\|_1)^2} \\ &\geq \frac{f(\mathbf{x}^t) - f(\mathbf{x}^*)}{8n^2M^2(B + (3/2)B)^2} \\ &\geq \frac{f(\mathbf{x}^t) - f(\mathbf{x}^*)}{50n^2M^2B^2}, \end{aligned}$$

where we used the fact that  $\|\mathbf{x}^*\|_1 \leq n \|\mathbf{x}^*\|_\infty \leq nB$ . Equivalently,

$$f(\mathbf{x}^{t+1}) - f(\mathbf{x}^*) \leq \left( 1 - \frac{1}{50n^2M^2B^2} \right) (f(\mathbf{x}^t) - f(\mathbf{x}^*)),$$

and summing up these for  $t \in \{0, 1, \dots, \bar{t} - 1\}$ , we get

$$\begin{aligned} f(\mathbf{x}^{\bar{t}}) - f(\mathbf{x}^*) &\leq \left(1 - \frac{1}{50n^2M^2B^2}\right)^{\bar{t}} (f(\mathbf{x}^0) - f(\mathbf{x}^*)) \\ &\leq \varepsilon, \end{aligned}$$

as long as

$$\bar{t} \geq 50n^2M^2B^2 \log \frac{f(\mathbf{x}^0) - f(\mathbf{x}^*)}{\varepsilon},$$

therefore we conclude that  $\bar{t}$  is at most this quantity.

Now we consider the iterations  $t \geq \bar{t}$ . If  $f(\mathbf{x}^t) \leq f(\mathbf{x}^*) + \varepsilon$  there are no such iterations and we are done. Therefore we have that  $f(\mathbf{x}^t) \leq 2f(\mathbf{x}^*)$ . We again use Lemma A.2 for all such  $t$ , which gives

$$\begin{aligned} f(\mathbf{x}^t) - f(\mathbf{x}^{t+1}) &\geq \frac{(f(\mathbf{x}^t) - f(\mathbf{x}^*))^2}{4M^2f(\mathbf{x}^t)(\|\mathbf{x}^*\|_1 + \|\mathbf{x}^t\|_1)^2} \\ &\geq \frac{(f(\mathbf{x}^t) - f(\mathbf{x}^*))^2}{25f(\mathbf{x}^t)n^2M^2B^2} \\ &\geq \frac{(f(\mathbf{x}^t) - f(\mathbf{x}^*))^2}{50f(\mathbf{x}^*)n^2M^2B^2}. \end{aligned}$$

By known convergence results, this recurrence leads to the bound

$$\begin{aligned} f(\mathbf{x}^T) &\leq f(\mathbf{x}^*) + \frac{100f(\mathbf{x}^*)n^2M^2B^2}{T - \bar{t}} \\ &\leq f(\mathbf{x}^*) \left(1 + \frac{100n^2M^2B^2}{T - \bar{t}}\right), \end{aligned}$$

implying that  $f(\mathbf{x}^T) \leq (1 + \delta)f(\mathbf{x}^*)$  after

$$T - \bar{t} = O\left(n^2M^2B^2\frac{1}{\delta}\right)$$

additional iterations after  $\bar{t}$ . Therefore, the total number of iterations to achieve  $f(\mathbf{x}^T) \leq (1 + \delta) \cdot f(\mathbf{x}^*) + \varepsilon$  is

$$O\left(n^2M^2B^2\left(\frac{1}{\delta} + \log \frac{f(\mathbf{x}^0) - f(\mathbf{x}^*)}{\varepsilon}\right)\right).$$

□

## A.2.2 PROOF OF THEOREM 4.1

*Proof.* Similarly to the proof of Corollary 5.1, we apply Lemma A.2 for  $T$  iterations to obtain solutions  $\mathbf{x}^0, \dots, \mathbf{x}^T$ , but now we also have to account for the sparsity increase of  $\mathbf{x}^t$ . For this reason, we use  $\lambda_t < 1$ , which disincentivizes updating zero entries of the solution vector.

Compared to Corollary 5.1, we have the differences that

$$\lambda_t^{-1} \|\mathbf{x}^*\|_1 = \max\{c^{-1} \|\mathbf{x}^t\|_1, \|\mathbf{x}^*\|_1\},$$

and that we have the following tighter bounds because of sparsity:

$$\begin{aligned} \|\mathbf{x}^*\|_1 &\leq sB \\ \|\mathbf{x}^t\|_1 &\leq \|\mathbf{x}^t\|_0 \|\mathbf{x}^t\|_\infty \leq \|\mathbf{x}^t\|_0 (3/2)B. \end{aligned}$$

We first bound the sparsity. Note that the sparsity increases by at most 1 every time the first bullet of Lemma A.2 is true, and does not increase when the second bullet is true. Therefore, the progress in each sparsity-increasing iteration is

$$\begin{aligned} f(\mathbf{x}^t) - f(\mathbf{x}^{t+1}) &\geq \frac{(f(\mathbf{x}^t) - f(\mathbf{x}^*))^2}{4f(\mathbf{x}^*)M^2(\|\mathbf{x}^*\|_1 + \lambda_t \|\mathbf{x}^t\|_1)^2} \\ &\geq \frac{(f(\mathbf{x}^t) - f(\mathbf{x}^*))^2}{4(1+c)^2f(\mathbf{x}^*)M^2\|\mathbf{x}^*\|_1^2}. \end{aligned}$$

Completely analogously to the proof of Corollary 5.1, this implies that the total number of such iterations (and thus total sparsity) is

$$s' := \|\mathbf{x}^T\|_0 \leq O\left(\|\mathbf{x}^*\|_1^2 (1+c)^2 M^2 \left(\frac{1}{\delta} + \log \frac{f(\mathbf{x}^0) - f(\mathbf{x}^*)}{\varepsilon}\right)\right).$$

Now it remains to bound the total number of iterations. We have

$$\begin{aligned} & \max\{\|\mathbf{x}^*\|_1 + \lambda_t \|\mathbf{x}^t\|_1, \lambda_t^{-1} \|\mathbf{x}^*\|_1 + \|\mathbf{x}^t\|_1\} \\ & \leq \max\{\|\mathbf{x}^*\|_1 + \|\mathbf{x}^t\|_1, c^{-1} \|\mathbf{x}^t\|_1 + \|\mathbf{x}^t\|_1\} \\ & \leq \|\mathbf{x}^*\|_1 + (1+c^{-1}) \|\mathbf{x}^t\|_1 \\ & \leq \|\mathbf{x}^*\|_1 + \frac{3}{2}(1+c^{-1}) \|\mathbf{x}^t\|_0 B \\ & \leq \|\mathbf{x}^*\|_1 + \frac{3}{2}(1+c^{-1}) \|\mathbf{x}^T\|_0 B. \end{aligned}$$

As a result, the progress bound of Lemma A.2 becomes

$$f(\mathbf{x}^t) - f(\mathbf{x}^{t+1}) \geq \frac{(f(\mathbf{x}) - f(\mathbf{x}^*))^2}{4f(\mathbf{x})M^2 \left(\|\mathbf{x}^*\|_1 + \frac{3}{2}(1+c^{-1}) \|\mathbf{x}^T\|_0 B\right)^2},$$

and, analogously to the proof of Corollary 5.1 and using the fact that  $\|\mathbf{x}^*\|_1 \leq \|\mathbf{x}^*\|_0 B$ , the total number of iterations is bounded by

$$T = O\left((1+c^{-1})^2 \left(\|\mathbf{x}^T\|_0^2 + \|\mathbf{x}^*\|_0^2\right) M^2 B^2 \left(\frac{1}{\delta} + \log \frac{f(\mathbf{x}^0) - f(\mathbf{x}^*)}{\varepsilon}\right)\right).$$

□

### A.2.3 PROOF OF THEOREM 4.3

*Proof.* We move similarly to the proof of Theorem 4.1, but now we can strengthen Lemma A.2 because  $\mathbf{x}^t$  is fully corrected for all  $t$ , i.e.  $\nabla_i f(\mathbf{x}^t) = 0$  for all  $i \in \text{supp}(\mathbf{x}^t)$ . As in the proof of Lemma A.2, we can lower bound the amount of progress as a function of  $\|\nabla f(\mathbf{x}^t)\|_\infty$  as follows:

$$f(\mathbf{x}^t) - f(\mathbf{x}^{t+1}) \geq \min\left\{\frac{(\nabla_i f(\mathbf{x}^t))^2}{4M^2 f(\mathbf{x}^t)}, \frac{|\nabla_i f(\mathbf{x}^t)|}{4M}\right\}.$$

Now, by convexity of  $f$  we have

$$\langle \nabla f(\mathbf{x}^t), \mathbf{x}^t - \mathbf{x}^* \rangle \geq f(\mathbf{x}^t) - f(\mathbf{x}^*). \quad (7)$$

Because of fully corrective steps we have  $\langle \nabla f(\mathbf{x}^t), \mathbf{x}^t \rangle = 0$ , and so the left hand side of (7) is upper bounded by  $\|\nabla f(\mathbf{x}^t)\|_\infty \|\mathbf{x}^*\|_1$ . As a result, we have

$$\|\nabla f(\mathbf{x}^t)\|_\infty^2 \geq \frac{(f(\mathbf{x}) - f(\mathbf{x}^*))^2}{\|\mathbf{x}^*\|_1^2},$$

and so we get the progress bound of

$$\begin{aligned} f(\mathbf{x}^t) - f(\mathbf{x}^{t+1}) & \geq \min\left\{\frac{(f(\mathbf{x}^t) - f(\mathbf{x}^*))^2}{4M^2 f(\mathbf{x}^t) \|\mathbf{x}^*\|_1^2}, \frac{f(\mathbf{x}^t) - f(\mathbf{x}^*)}{4M \|\mathbf{x}^*\|_1}\right\} \\ & \geq \frac{(f(\mathbf{x}^t) - f(\mathbf{x}^*))^2}{4M^2 f(\mathbf{x}^t) \|\mathbf{x}^*\|_1^2}, \end{aligned}$$

because  $M \|\mathbf{x}^*\|_1 > 1$ . Similarly to the proof of Theorem 4.1, this progress bound leads to a sparsity of

$$s' := \|\mathbf{x}^T\|_0 \leq O\left(\|\mathbf{x}^*\|_1^2 M^2 \left(\frac{1}{\delta} + \log \frac{f(\mathbf{x}^0) - f(\mathbf{x}^*)}{\varepsilon}\right)\right)$$

and the same number of iterations. □

## B MISSING PROOFS FROM SECTION 5

### B.1 GRADIENT UPDATE LEMMA

**Lemma B.1** (Gradient update). *Let  $f : \mathbb{R}^n \rightarrow \mathbb{R}_{>0}$  be a twice continuously differentiable convex function that is  $2\gamma$ -second order robust with respect to the  $\ell_1$  norm and  $\mu$ -multiplicatively smooth with respect to the  $\ell_2$  norm for some  $\gamma, \mu > 0$ . Let  $\mathbf{x} \in \mathbb{R}^n$  be a solution such that  $f(\mathbf{x}) > f(\mathbf{x}^*)$ , where  $\mathbf{x}^* \in \mathbb{R}^n$  is an unknown solution with  $\|\mathbf{x} - \mathbf{x}^*\|_2 \leq \|\mathbf{x}^*\|_2$ . We make the update*

$$\mathbf{x}' = \mathbf{x} - \eta \nabla f(\mathbf{x}),$$

where  $\eta = 0.5 \min \left\{ \frac{1}{\mu f(\mathbf{x})}, \frac{1}{\gamma \|\nabla f(\mathbf{x})\|_1} \right\}$  is a step size. Then, the progress in decreasing  $f$  is:

$$f(\mathbf{x}) - f(\mathbf{x}') \geq \min \left\{ \frac{(f(\mathbf{x}) - f(\mathbf{x}^*))^2}{4\mu f(\mathbf{x}) \|\mathbf{x}^*\|_2^2}, \frac{f(\mathbf{x}) - f(\mathbf{x}^*)}{4\gamma \sqrt{n} \|\mathbf{x}^*\|_2} \right\}.$$

Additionally, as long as  $\mathbf{x}'$  is still suboptimal with respect to  $\mathbf{x}^*$ , i.e.  $f(\mathbf{x}') > f(\mathbf{x}^*)$ , the distance to  $\mathbf{x}^*$  decreases:  $\|\mathbf{x}' - \mathbf{x}^*\|_2 \leq \|\mathbf{x} - \mathbf{x}^*\|_2$ . Finally, if  $f(\mathbf{x}) \leq f(\mathbf{x}^*)$ , then  $f(\mathbf{x}') \leq f(\mathbf{x})$ .

*Proof.* We first consider a generic update  $\mathbf{x}' = \mathbf{x} + \tilde{\mathbf{x}}$ . By Taylor's theorem and the fact that  $f$  is twice continuously differentiable, we have

$$f(\mathbf{x}') = f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \tilde{\mathbf{x}} \rangle + \frac{1}{2} \langle \tilde{\mathbf{x}}, \nabla^2 f(\bar{\mathbf{x}}) \tilde{\mathbf{x}} \rangle,$$

for some  $\bar{\mathbf{x}}$  that is entrywise between  $\mathbf{x}$  and  $\mathbf{x}'$ .

Since  $f$  is  $2\gamma$ -second-order-robust with respect to  $\ell_1$  and  $\mu$ -multiplicatively-smooth with respect to the  $\ell_2$  norm, as long as the update is bounded in  $\ell_1$  norm as

$$\|\tilde{\mathbf{x}}\|_1 \leq 1/(2\gamma) \tag{8}$$

we have

$$\begin{aligned} f(\mathbf{x}') &\leq f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \tilde{\mathbf{x}} \rangle + \langle \tilde{\mathbf{x}}, \nabla^2 f(\mathbf{x}) \tilde{\mathbf{x}} \rangle \\ &\leq f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \tilde{\mathbf{x}} \rangle + \mu f(\mathbf{x}) \|\tilde{\mathbf{x}}\|_2^2. \end{aligned}$$

Note that the right hand side is minimized for

$$\tilde{\mathbf{x}} = -\frac{1}{2\mu f(\mathbf{x})} \nabla f(\mathbf{x}).$$

In addition, to stay in the  $\ell_1$  neighborhood where the Hessian is stable, we need to satisfy (8). Based on these requirements, we make the update  $\tilde{\mathbf{x}} = -\eta \nabla f(\mathbf{x})$ , where

$$\eta = \min \left\{ \frac{1}{2\mu f(\mathbf{x})}, \frac{1}{2\gamma \|\nabla f(\mathbf{x})\|_1} \right\}.$$

We thus have

$$\begin{aligned} f(\mathbf{x}) - f(\mathbf{x}') &\geq (\eta - \eta^2 \mu f(\mathbf{x})) \|\nabla f(\mathbf{x})\|_2^2 \\ &\geq \frac{\eta}{2} \|\nabla f(\mathbf{x})\|_2^2 \end{aligned}$$

and so

$$\begin{aligned} f(\mathbf{x}) - f(\mathbf{x}') &\geq \min \left\{ \frac{1}{4\mu f(\mathbf{x})}, \frac{1}{4\gamma \|\nabla f(\mathbf{x})\|_1} \right\} \|\nabla f(\mathbf{x})\|_2^2 \\ &\geq \min \left\{ \frac{\|\nabla f(\mathbf{x})\|_2^2}{4\mu f(\mathbf{x})}, \frac{\|\nabla f(\mathbf{x})\|_2}{4\gamma \sqrt{n}} \right\}. \end{aligned}$$



This takes care of the case  $f(\mathbf{x}) \leq f(\mathbf{x}^*)$ , since it shows that  $f(\mathbf{x}') \leq f(\mathbf{x})$ . Now we deal with the case  $f(\mathbf{x}) > f(\mathbf{x}^*)$ . By convexity we have

$$\begin{aligned} f(\mathbf{x}^*) &\geq f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \mathbf{x}^* - \mathbf{x} \rangle \\ &\geq f(\mathbf{x}) - \|\nabla f(\mathbf{x})\|_2 \|\mathbf{x}^* - \mathbf{x}\|_2 \\ &\geq f(\mathbf{x}) - \|\nabla f(\mathbf{x})\|_2 \|\mathbf{x}^*\|_2, \end{aligned}$$

which gives

$$\|\nabla f(\mathbf{x})\|_2^2 \geq \frac{(f(\mathbf{x}) - f(\mathbf{x}^*))^2}{\|\mathbf{x}^*\|_2^2},$$

and so

$$f(\mathbf{x}) - f(\mathbf{x}') \geq \min \left\{ \frac{(f(\mathbf{x}) - f(\mathbf{x}^*))^2}{4\mu f(\mathbf{x}) \|\mathbf{x}^*\|_2^2}, \frac{f(\mathbf{x}) - f(\mathbf{x}^*)}{4\gamma\sqrt{n} \|\mathbf{x}^*\|_2} \right\}.$$

For the norm bound, we suppose that  $f(\mathbf{x}') > f(\mathbf{x}^*)$  (otherwise we are done). We have

$$\begin{aligned} &\|\mathbf{x}' - \mathbf{x}^*\|_2^2 - \|\mathbf{x} - \mathbf{x}^*\|_2^2 \\ &= \|\mathbf{x}' - \mathbf{x}\|_2^2 + 2\langle \mathbf{x} - \mathbf{x}^*, \mathbf{x}' - \mathbf{x} \rangle \\ &= \eta^2 \|\nabla f(\mathbf{x})\|_2^2 - 2\eta \langle \mathbf{x} - \mathbf{x}^*, \nabla f(\mathbf{x}) \rangle. \end{aligned}$$

Now, note that

$$\frac{\eta}{2} \|\nabla f(\mathbf{x})\|_2^2 \leq f(\mathbf{x}) - f(\mathbf{x}') \leq f(\mathbf{x}) - f(\mathbf{x}^*)$$

and by convexity  $\langle \mathbf{x} - \mathbf{x}^*, \nabla f(\mathbf{x}) \rangle \geq f(\mathbf{x}) - f(\mathbf{x}^*)$ , so

$$\begin{aligned} &\|\mathbf{x}' - \mathbf{x}^*\|_2^2 - \|\mathbf{x} - \mathbf{x}^*\|_2^2 \\ &= \eta^2 \|\nabla f(\mathbf{x})\|_2^2 - 2\eta \langle \mathbf{x} - \mathbf{x}^*, \nabla f(\mathbf{x}) \rangle \\ &\leq 0. \end{aligned}$$

□

## B.2 PROOF OF THEOREM 5.2

*Proof.* We repeatedly use Lemma B.1 to obtain iterates  $\mathbf{x}^0, \mathbf{x}^1, \dots, \mathbf{x}^T$ . Note that as long as  $f(\mathbf{x}^t) > f(\mathbf{x}^*)$ , we have  $\|\mathbf{x}^t - \mathbf{x}^*\|_2 \leq \|\mathbf{x}^0 - \mathbf{x}^*\|_2 := R$ .

Let  $\bar{t}$  be the smallest  $t \geq 0$  for which  $f(\mathbf{x}^{\bar{t}}) \leq 2f(\mathbf{x}^*)$  or  $f(\mathbf{x}^{\bar{t}}) \leq f(\mathbf{x}^*) + \varepsilon$ , and let  $\bar{t} = \infty$  if this never happens. Therefore, for all  $t < \bar{t}$  we have  $f(\mathbf{x}^t) \geq 2f(\mathbf{x}^*) \Rightarrow \frac{f(\mathbf{x}^t) - f(\mathbf{x}^*)}{f(\mathbf{x}^t)} \geq \frac{1}{2}$ , and so the statement of Lemma B.1 gives:

$$\begin{aligned} f(\mathbf{x}^t) - f(\mathbf{x}^{t+1}) &\geq \min \left\{ \frac{1}{8\mu \|\mathbf{x}^*\|_2^2}, \frac{1}{4\gamma\sqrt{n} \|\mathbf{x}^*\|_2} \right\} \cdot (f(\mathbf{x}^t) - f(\mathbf{x}^*)) \\ &\geq \frac{1}{8\mu \|\mathbf{x}^*\|_2^2 + 4\gamma\sqrt{n} \|\mathbf{x}^*\|_2} \cdot (f(\mathbf{x}^t) - f(\mathbf{x}^*)). \end{aligned}$$

Equivalently,

$$f(\mathbf{x}^{t+1}) - f(\mathbf{x}^*) \leq \left( 1 - \frac{1}{8\mu \|\mathbf{x}^*\|_2^2 + 4\gamma\sqrt{n} \|\mathbf{x}^*\|_2} \right) (f(\mathbf{x}^t) - f(\mathbf{x}^*)),$$

and summing up these for  $t \in \{0, 1, \dots, \bar{t} - 1\}$ , we get

$$\begin{aligned} f(\mathbf{x}^{\bar{t}}) - f(\mathbf{x}^*) &\leq \left( 1 - \frac{1}{8\mu \|\mathbf{x}^*\|_2^2 + 4\gamma\sqrt{n} \|\mathbf{x}^*\|_2} \right)^{\bar{t}} (f(\mathbf{x}^0) - f(\mathbf{x}^*)) \\ &\leq \varepsilon, \end{aligned}$$

as long as

$$\bar{t} \geq (8\mu R^2 + 4\gamma\sqrt{n}R) \log \frac{f(\mathbf{x}^0) - f(\mathbf{x}^*)}{\varepsilon},$$

therefore we conclude that  $\bar{t}$  is at most this quantity.

Now we consider the iterations  $t \geq \bar{t}$ . If  $f(\mathbf{x}^{\bar{t}}) \leq f(\mathbf{x}^*) + \varepsilon$  there are no such iterations and we are done. Therefore we have that  $f(\mathbf{x}^{\bar{t}}) \leq 2f(\mathbf{x}^*)$ . We again use Lemma B.1 for all such  $t$ , which gives

$$\begin{aligned} f(\mathbf{x}^t) - f(\mathbf{x}^{t+1}) &\geq \frac{(f(\mathbf{x}^t) - f(\mathbf{x}^*))^2}{4\mu f(\mathbf{x}^t)R^2} \\ &\geq \frac{(f(\mathbf{x}^t) - f(\mathbf{x}^*))^2}{8\mu f(\mathbf{x}^*)R^2}. \end{aligned}$$

By known convergence results, this recurrence leads to the bound

$$\begin{aligned} f(\mathbf{x}^T) &\leq f(\mathbf{x}^*) + \frac{16\mu f(\mathbf{x}^*)R^2}{T - \bar{t}} \\ &= f(\mathbf{x}^*) \left(1 + \frac{16\mu R^2}{T - \bar{t}}\right), \end{aligned}$$

implying that  $f(\mathbf{x}^T) \leq (1 + \delta)f(\mathbf{x}^*)$  after

$$T - \bar{t} = O\left(\mu R^2 \frac{1}{\delta}\right)$$

additional iterations after  $\bar{t}$ .

Therefore, the total number of iterations to achieve  $f(\mathbf{x}^T) \leq (1 + \delta) \cdot f(\mathbf{x}^*) + \varepsilon$  is

$$O\left((\mu R^2 + \gamma\sqrt{n}R) \left(\frac{1}{\delta} + \log \frac{f(\mathbf{x}^0) - f(\mathbf{x}^*)}{\varepsilon}\right)\right).$$

For  $\mu = \beta m^{-1}$ ,  $\gamma = \sqrt{\beta}$ , and using the fact that  $R \geq \sqrt{n}$ , we get

$$O\left(\frac{\beta R^2}{m} \left(\frac{1}{\delta} + \log \frac{f(\mathbf{x}^0) - f(\mathbf{x}^*)}{\varepsilon}\right)\right)$$

iterations. □

## C PROOF OF THEOREM 6.1

*Proof.* Let us consider a classifier  $\mathbf{x}^*$  with  $\|\mathbf{x}^*\|_2 = 1$  and margins  $\geq \alpha$ , i.e.  $b_i(\mathbf{A}\mathbf{x}^*)_i \geq \alpha$  for all  $i \in [m]$ . Now, Corollary 5.1 and Theorem 5.2 imply that we can compute a solution  $f_\lambda(\mathbf{x}) \leq 2f_\lambda(\mathbf{x}^*) + \varepsilon$  after  $T = O(\lambda^2 X \log \frac{m}{\varepsilon})$  iterations. Now, note that

$$\begin{aligned} \sum_i \log(1 + e^{-\lambda b_i(\mathbf{A}\mathbf{x})_i}) &\leq 2 \sum_i \log(1 + e^{-\lambda b_i(\mathbf{A}\mathbf{x}^*)_i}) + \varepsilon \\ &\leq 2m \log(1 + e^{-\lambda\alpha}) + \varepsilon \\ &\leq 2me^{-\lambda\alpha} + \varepsilon \\ &\leq 3me^{-\lambda\alpha}, \end{aligned}$$

after setting  $\varepsilon = me^{-\lambda\alpha}$ .

Now, re-arranging and using the fact that  $3me^{-\lambda\alpha} \leq 2$  implies  $e^{3me^{-\lambda\alpha}} \leq 1 + 6me^{-\lambda\alpha}$ , we have that

$$\begin{aligned} b_i(\mathbf{A}\mathbf{x})_i &\geq -\lambda^{-1} \log(e^{3me^{-\lambda\alpha}} - 1) \\ &\geq -\lambda^{-1} \log(6me^{-\lambda\alpha}) \\ &= \alpha - \lambda^{-1} \log(6m) \\ &\geq \alpha(1 - \hat{\varepsilon}), \end{aligned}$$

where the last equality follows by our setting of  $\lambda \geq \frac{\log(6m)}{\alpha\hat{\varepsilon}}$ . Therefore, the number of iterations is  $O(\lambda^3\alpha) \leq \tilde{O}\left(\frac{1}{\alpha^2\hat{\varepsilon}^3}\right)$ . Additionally,  $\|\mathbf{x}\|_2 \leq \|\mathbf{x}^*\|_2$ .

To bound the distance from the classifier, we note that

$$E^2 := \left\| \frac{\mathbf{x}}{\|\mathbf{x}\|_2} - \mathbf{x}^* \right\|_2^2 = \left\| \frac{\mathbf{x}}{\|\mathbf{x}\|_2} - \mathbf{x}^* \right\|_2^2 = 2 - 2 \left\langle \frac{\mathbf{x}}{\|\mathbf{x}\|_2}, \mathbf{x}^* \right\rangle.$$

On the other hand, we let  $\bar{\mathbf{x}} = \frac{1}{2} \frac{\mathbf{x}}{\|\mathbf{x}\|_2} + \frac{1}{2} \mathbf{x}^*$  and compute its smallest margin as

$$\alpha \geq \frac{b_i(\mathbf{A}\bar{\mathbf{x}})_i}{\|\bar{\mathbf{x}}\|_2} \geq \frac{\frac{\alpha(1-\hat{\varepsilon})}{2\|\mathbf{x}\|_2} + \frac{\alpha}{2}}{\sqrt{\frac{1}{4} + \frac{1}{4} + \frac{1}{2} \left\langle \frac{\mathbf{x}}{\|\mathbf{x}\|_2}, \mathbf{x}^* \right\rangle}} \geq \frac{\frac{\alpha(1-\hat{\varepsilon})}{2} + \frac{\alpha}{2}}{\sqrt{1 - E^2/4}}.$$

Re-arranging, we get that

$$\begin{aligned} 1 - E^2/4 &\geq \frac{1}{4}(2 - \hat{\varepsilon})^2 = 1 - \hat{\varepsilon} + \hat{\varepsilon}^2/4 \\ E^2 &\leq 4\hat{\varepsilon}(1 - \hat{\varepsilon}) \leq 4\hat{\varepsilon}. \end{aligned}$$

Therefore, we have

$$\left\| \frac{\mathbf{x}}{\|\mathbf{x}\|_2} - \mathbf{x}^* \right\|_2 \leq 2\sqrt{\hat{\varepsilon}},$$

or in other words

$$\left\| \frac{\mathbf{x}}{\|\mathbf{x}\|_2} - \mathbf{x}^* \right\|_2 \leq E$$

after  $\tilde{O}\left(\frac{1}{\alpha^2 E^6}\right)$  iterations. □