

1 Before we begin, we highlight one of our key contributions: **Algorithm 2**. It showcases a remarkable *interplay between*  
 2 *statistics and optimization*: the *increasing* step sizes scheme (required for computational optimality) only works because  
 3 we rely on early stopping and do not aim to fully optimize the training objective. In contrast, most results in optimization  
 4 literature consider gradient descent with constant or *decreasing* step sizes to ensure convergence to the objective.

5 **Providing more intuition (R#1 and R#2)**. For gradient descent (GD), our reparameterization turns *additive* updates  
 6 into *multiplicative* updates (see lines 69-74). As a result, the scale of the parameter can be understood as inertia – the  
 7 small parameters have a tendency to remain almost unchanged, while the larger parameters are more sensitive to the  
 8 gradient size with respect to the standard parameterization  $\mathbf{w}$ . Sparsity is induced with reparameterization *together* with  
 9 small initialization size (one without the other doesn't work). For more intuition see the proof sketch (Theorem 4),  
 10 simulations and appendices A and B. Finally, previous work in the literature shows that GD implicitly regularizes the  $\ell_2$   
 11 norm. This corresponds to minimizing  $\ell_1$  norm of  $\mathbf{w}$  in our parameterization on  $\mathbf{u}, \mathbf{v}$  (see lines 255-256).

12 **Choice of hyperparameters (R#1 and R#2)**. As discussed in lines 201-208, we only need to know  $w_{\max}^*$  up to  
 13 multiplicative factors to properly initialize  $\alpha$  and  $\eta$ . Theorem 2 shows how to obtain such an estimate. Hence we only  
 14 need to tune the stopping time, which can be done by cross-validation.

15 **Response to Reviewer #1.**

16 **1.** See our paragraphs on intuition and hyperparameters above.

17 **2.** Most work on implicit  $\ell_2$  regularization focus on GD with a constant step size usually stopping at  $\Theta(\sqrt{n})$  iterations;  
 18  $(\eta t)^{-1}$  corresponds to the Ridge regression  $\lambda$ . In lines 195-200 we discuss connections to Thm 1 and 3. On the other  
 19 hand, we are not aware of other work on implicit regularization achieving computational optimality via an increasing  
 20 step sizes scheme (Alg. 2). We highlight that in our case implicit regularizer is not strictly  $\ell_1$  norm (see lines 125-126  
 21 and 334-339) and our work, to the best of our knowledge, is first to induce sparsity *implicitly* in a general noisy setting.

22 **3.** This remains an open question not considered in our paper. We believe that a good starting point would be to  
 23 experiment with individualized initialization sizes and step sizes among each dimension/group.

24 **Improvements section.** We agree with the suggestion and plan to add an additional paragraph to the related literature  
 25 section, expanding on the second point above. Subject to space considerations we will also expand on intuition.

26 **Response to Reviewer #2.**

27 **1.** For sparsity of the optimization path see proof sketch (Thm 4), simulations (lines 314-316)  
 28 and the main proofs. For sparsity at the stopping time, see the  $\ell_\infty$  bound on  $S^c$  in Thm 1.

29 **2.**  $w_{\max}^*$  is defined in line 84. Line 151 refers to table of notation.

30 **3.** Since  $\mathbf{X}$  needs to only satisfy RIP,  $n$  depends only logarithmically on  $d$ .

31 **Improvements section.** For intuition and hyperparameters, see the two paragraphs at the  
 32 top of the rebuttal. To address the concerns on RIP assumption being too strong, we have  
 33 performed additional situations when RIP assumption fails. Consider the setting given in  
 34 lines 322-332, with rows of  $\mathbf{X}$  now sampled from  $N(\mathbf{0}, \Sigma)$ , with  $\Sigma = (1 - \mu)\mathbf{I} + \mu\mathbf{1}\mathbf{1}^\top/d$ .  
 35 On the right, we plot simulation results with  $\mu = 0$  (RIP holds) and  $\mu = 0.5$  (RIP fails). We  
 36 see that even when RIP fails, our method still exhibits correct rates and outperforms the lasso  
 37 when the phase transition happens. The gap between gradient descent and the oracle method  
 38 is visible due to the  $\log k$  factor in Corollary 3, suggesting also that the rate given in Corollary 3 could be tight. We  
 39 will address the reviewers concerns by adding a section on potential improvements with an expansion of the above  
 40 discussion. We will also compare and contrast RIP and RE assumptions. If space permits, we will also slightly expand  
 41 on the intuition.

42 **Response to Reviewer #3.**

43 We have previously attributed the quadratic sample complexity in  $k$  to our bounds  
 44 being  $\ell_\infty$  (which is harder) rather than  $\ell_2$ . Our focus has been on *minimax-rates and*  
 45 *dimension-independent rates* with *optimal computational complexity*. Also, while there  
 46 is loss in sample complexity, there is gain in performance that is impossible to be  
 47 achieved by the lasso (see Corollary 3 and lines 334-339).

48 We stick to the simulation setting described in lines 308-313 and 322-328, with  $d =$   
 49  $5000$ . The left figure on the right compares  $\ell_2$  error ratios for gradient descent and  
 50 lasso. The blue region corresponds to our method achieving lower error, while the red  
 51 region corresponds to the lasso achieving lower error than gradient descent. This plot strongly suggests, that sample  
 52 complexity linear in  $k$  should indeed be enough to match/exceed performance of the lasso. The question remains,  
 53 whether the  $\ell_\infty$  bounds in Theorems 1 and 3 (in particular for stopping time  $t$ ,  $\|\mathbf{w}_t \odot \mathbf{1}_{S^c}\|_\infty \leq \sqrt{\alpha}$ ) require sample  
 54 complexity quadratic in  $k$ ? The figure on the right side suggests that the sample complexity linear in  $k$  is enough to  
 55 satisfy even the  $\ell_\infty$  bounds. We expect this sample complexity gap to be addressed in future work.

56 Given the results above, we absolutely agree with the suggestion to include a discussion on sub-optimal sample  
 57 complexity in our revision. We plan to do so in an extra section on potential improvements (see also response to R#2).

