

1 We first and foremost thank the reviewers for their valuable time and feedback. We will fix all small comments and  
2 typos the reviewers note.

3 Reviewer #1 asks how our results relate to adaptive gradient methods. As the reviewer notes, adaptive gradient methods  
4 construct a sequence of linear pre-conditioners to be applied to the stochastic sub-gradient before the update. Our  
5 results certify that there exists an optimal linear pre-conditioner for quadratically convex constraint sets. As such,  
6 adaptive gradient methods can be minimax (rate) optimal. In online algorithms, the common practice [4, 5, 6, 7, 2]  
7 is to measure regret with respect to the “best” post-hoc regularizer (i.e. preconditioner) and fixed predictor; in this  
8 context, the regret achieved by AdaGrad is a factor  $\sqrt{2}$  away from the regret the best post-hoc linear pre-conditioner  
9 achieves over rectangular domains (and may be  $\sqrt{d}$  better than standard gradient methods). Our results thus guarantee  
10 the minimax optimality of AdaGrad in certain settings. We are perhaps a bit imprecise in that we use “best linear  
11 preconditioner” as a shorthand for what adaptive algorithms may achieve; we will make this more precise. The extent  
12 to which specific adaptive algorithms find the (optimal) linear pre-conditioner for specific constraint sets remains open.

13 Reviewer #1 correctly notes that the quadratic convexity of the constraint set is critical via Proposition 4. In the case  
14 that  $\Theta$  is not quadratically convex, one must replace  $\Theta$  by its quadratic hull when swapping the infimum and supremum,  
15 resulting in a (potentially) much larger set (e.g.  $\text{QHull}(\mathbf{B}_1(0, 1)) = \mathbf{B}_2(0, 1)$ ). We will emphasize this.

16 Reviewer #1 asks about the origin and meaning of Corollary 3. It follows from Corollary 2 by lower bounding the  
17 inequality. To illustrate the corollary, we can observe that when  $\gamma$  is an  $\ell_q$  norm with  $q \in [1, 2]$ , the lower and upper  
18 bounds match up to  $\sqrt{\log d}$ . After the submission of this paper, we derived more precise bounds in the case where  
19 the  $\gamma$ -ball is not quadratically convex: we obtained matching lower and upper bounds when the gradients live in any  
20 weighted  $\ell_q$  ball i.e.  $\gamma(g) = (\sum_{j \leq d} a_j |g_j|^q)^{1/q}$  for  $a_j > 0$  and  $q \in [1, \infty]$ . We will include these new results.

21 Reviewer #1 asks for definitional clarifications for minimax risk and regret. While  $F_P$  is a deterministic function, in  
22 the definition of the minimax risk, it is applied to  $\hat{\theta}_n(X_1^n)$  which is a (random) estimator based on a sample  $X_1^n \sim P$ ,  
23 necessitating an expectation. In the definition of  $\mathfrak{M}_n^S$  and  $\mathfrak{M}_n^R$ , the supremum over the sample space takes place before  
24 the infimum as the infimum ranges over all (measurable) functions  $\mathcal{X}^n \rightarrow \Theta$ . This definition accords with the literature  
25 on lower bounds in convex optimization, where the supremum is over stochastic oracles [1]. Reviewer #2 asks for  
26 a clarification on the definition of  $X_1^n$  and  $x_1^n$ . The former corresponds to a collection of  $n$  i.i.d. random variables  
27  $(X_1, \dots, X_n)$ . The latter denotes  $n$  (fixed) vectors  $(x_1, \dots, x_n) \in \mathcal{X}^n$ . We will clarify all of these definitions.

28 Reviewer #2 asks for applicability for the non-convex setting. Empirically, even in non-convex settings, AdaGrad  
29 tends to outperform vanilla gradient methods when data is sparse (e.g. [3, 8]). Our mathematical results probably do not  
30 translate as is beyond convexity. However, deriving similar results in the case, for example, of finding stationary points  
31 of non-convex functions is a natural extension and a very interesting future direction.

32 Reviewer #3 asks for concrete examples where the geometry of the constraint sets matters. In addition to the two deep  
33 learning examples of the above paragraph, this work is, for example, applicable in the broad case of linear models.  
34 In this setting, the constraint set corresponds to the set of classifiers of interest, and the geometry of the gradients  
35 corresponds to the geometry of the features (or covariates). For example, in NLP applications, bag-of-word features  
36 are very sparse by nature, so we seek a dense classifier (i.e. a weight for every word). In the terms of our paper, this  
37 means that  $\Theta$  is a weighted  $\ell_\infty$  ball,  $\gamma$  is a weighted  $\ell_1$  ball, and our theory suggests adaptive scaling is important. (For  
38 empirical results, see, e.g. [4].)

## 39 References

- 40 [1] A. Agarwal, P. L. Bartlett, P. Ravikumar, and M. J. Wainwright. Information-theoretic lower bounds on the oracle complexity of  
41 convex optimization. *IEEE Transactions on Information Theory*, 58(5):3235–3249, 2012.
- 42 [2] A. Cutcosky and T. Sarlos. Matrix-free preconditioning in online learning. In *Proceedings of the 36th International Conference*  
43 *on Machine Learning*, 2019.
- 44 [3] J. Dean, G. S. Corrado, R. Monga, K. Chen, M. Devin, Q. V. Le, M. Z. Mao, M. Ranzato, A. Senior, P. Tucker, K. Yang, and  
45 A. Y. Ng. Large scale distributed deep networks. In *Advances in Neural Information Processing Systems 26*, 2012.
- 46 [4] J. C. Duchi, E. Hazan, and Y. Singer. Adaptive subgradient methods for online learning and stochastic optimization. *Journal of*  
47 *Machine Learning Research*, 12:2121–2159, 2011.
- 48 [5] B. McMahan and M. Streeter. Adaptive bound optimization for online convex optimization. In *Proceedings of the Twenty Third*  
49 *Annual Conference on Computational Learning Theory*, 2010.
- 50 [6] F. Orabona and K. Crammer. New adaptive algorithms for online classification. In *Advances in Neural Information Processing*  
51 *Systems 23*, 2010.
- 52 [7] F. Orabona, K. Crammer, and N. Cesa-Bianchi. A generalized online mirror descent with applications to classification and  
53 regression. *Machine Learning*, 99(3):411–435, 2015.
- 54 [8] J. Pennington, R. Socher, and C. D. Manning. Glove: Global vectors for word representation. In *Proceedings of Empirical*  
55 *Methods for Natural Language Processing*, 2014.