

# 1 Summary

2 We appreciate the reviewers feedback! Generally, the reviewers suggestions could be decomposed into three categories:  
3 adding a related works section, cleaning up some of the notation, and clarifying and improving some of the examples  
4 and experiments. We address all three of these concerns below.

## 5 1.1 Adding a related works section

6 **Response** To address the concerns of reviewers 1 and 2, we will add a related work section to discuss how our frame-  
7 work relates to other acceleration frameworks proposed, including the ones listed by reviewer 2, and to clarify our con-  
8 tributions. In short, our work can be viewed as a kind of generalization of Allen-Zhu/Orecchia, Lessard/Recht/Packard,  
9 Lin/Mairal/Harchaoui, to more general  $p > 2$ , which allow us to obtain methods with faster convergence guarantees.  
10 Our work most closely resembles Wilson/Recht/Jordan, however we (1) introduce and discuss descent methods (which is  
11 omitted by Wilson/Recht/Jordan), and (2) provide a general description and Lyapunov analysis of the Monteiro-Svaiter  
12 acceleration framework. These manifold generalizations are what allow us to propose a novel method for optimization –  
13 RGD and ARGD – which has superior theoretical and empirical performance to several existing methods. We hope to  
14 make this clear in the related work section.

## 15 1.2 Clearing up notation

16 **Response** Reviewer 1: Equations (9) and (12) do indeed “hold” – perhaps the confusion is that we did not define  
17  $\|\cdot\|_{x_k}$ , which we will add. We set the  $B$  equal to the identity for the final three examples because for these examples,  
18 there is a natural definition of a norm given by the Hessian of the matrix  $h$ . We will change Lemma 4 to theorem 4. We  
19 will also add more details to the proofs so that it is easier to follow (although it would be helpful to understand which  
20 proofs the reviewer felt should be more detailed). Reviewer 3: We will add all notational suggestions. Thank you for  
21 the helpful clarifying suggestions.

## 22 1.3 Improving Examples and Experiments

23 **Response** Reviewer 1: We will clarify the notation DD and add a more detailed description of experimental results.  
24 Reviewer 2: The GLM loss (example 8) is a non-convex function. We are currently running experiments on MNIST for  
25 this objective and will add the results to the final version of our paper. As a preview, Hazan et al [Hazan et al., 2015, fig  
26 2] showed that the version of RGD (a.k.a stochastic normalized gradient descent) outperformed stochastic AGD on this  
27 objective. We hope to replicate this result in the deterministic setting to confirm our theoretical findings. Reviewer 3:  
28 You make an excellent point about the axis. Since we ran the experiment for a fixed  $10^{-6}$  iterations, we simply plotted  
29 the results without paying attention to whether it dropped below machine precision. We will cut off the plot at the stated  
30  $10^{-20}$ .

## 31 References

32 Elad Hazan, Kfir Y. Levy, and Shai Shalev-Shwartz. Beyond convexity: Stochastic quasi-convex optimization. In  
33 *Advances in Neural Information Processing Systems 28: Annual Conference on Neural Information Processing*  
34 *Systems 2015, December 7-12, 2015, Montreal, Quebec, Canada*, pages 1594–1602, 2015.