

We sincerely thank the four reviewers for their valuable feedback. We first discuss a few common concerns shared by reviewer 1, reviewer 2, reviewer 3 and reviewer 4.

• • **Non-convex bound:** Thanks for your constructive comments. It is clear that in convex case a better prediction reduces the bound. In the non-convex case it holds as well, with some careful analysis. For **H3**, if we alternatively consider $0 < m_t^T g_t = a \|g_t\|^2$ and $\|m_t\| \leq \|g_t\|$ (i.e. m_t lies in the hemisphere with g_t as its midline), we can show that \tilde{C}_2 reaches minimum when $a = 1$ (i.e. $m_t = g_t$). Also, \tilde{C}_1 is minimized at $a = 1$ under some conditions on the parameters (β_1, β_2 etc.). **That means the bound for non-convex case is tighter when m_t predicts g_t well, similar to the convex analysis.** We will adjust our discussion and presentation in the paper to address this point.

Reviewer 1: We thank the reviewer for valuable comments. Our point-to-point response is as follows:

Convex regret bound: For analysis purposes we presented the algorithm without projection step by assuming the compact assumption **H1**. Of course, this assumption needs to be verified and we partially did it for a model of interest that is a deep neural network, see Section 4.3. Adding projection steps is a neat idea to avoid having those issues but is not common in non-convex optimization analyses, see references [5, 9, 14, 38].

Numerical example: We thank the reviewer for their remark on the numerical runs. The main motivation behind those plots is to show that adding an optimistic update to the vanilla AMSGrad actually speed up the convergence in terms of both losses and accuracies. Given the well-known advantages of Adam-type methods as ADAM or AMSGrad, we did not compare to slower methods. We believe something like “SGD+optimistic” should be a future research worth studying. In this paper, we tried to focus on the most recent AMSGrad algorithm.

Reviewer 2: We thank the reviewer for valuable comments. A proofreading is being done we clarify that:

Novelty of the contribution: Although combining gradient prediction to AMSGrad update seems natural, as pointed out in the first paragraph of Section 3, we would like to stress on how the embedding of the prediction process (represented Figure 1) led to the two-stage algorithm OPT-AMSGrad (unlike the sequential structure of the original AMSGrad) where, first an auxiliary variable \tilde{w} is updated and then the global model w . Also, as discussed in the paper, optimistic learning is typically used in two-player-games, and to the best of our knowledge, this is the first proposal to apply optimistic acceleration to stochastic optimization problems (e.g. training deep neural networks).

Wall clock times: There are several works considering applying Alg. 3 in deep learning, e.g. [Nonlinear Acceleration of Deep Neural Networks, Scieur et al., 2018], with positive results. As noted in their paper, in practice extrapolation on CPU is faster than a forward pass on mini-batch and can be further accelerated on GPU. Moreover, note that at each iteration, we only change one past gradient, so we do not need to compute the whole linear system every time (but only a small portion). Therefore, it could be fairly efficient in practice. Secondly, the main focus of our paper is essentially the framework of integrating optimistic learning with AMSGrad. We chose Algorithm 3 mainly because of the empirical success reported in prior works. The choice of gradient prediction method is actually flexible. So, OPT-AMSGrad will definitely benefit from an algorithm with faster running time and good prediction quality. This is more related to acceleration literature and we regard this as an interesting line of future research.

Reviewer 3: We thank the reviewer for the thorough analysis. Our remarks are listed below:

Gradient prediction algorithm: We agree that a characterization on how well the gradient is predicted would be impactful. The scope of our paper being the *stochastic optimization problem* itself, we invoked a simple but effective gradient prediction algorithm on the basis of [31], which theory holds only for convex functions. As replied to Reviewer 2, we chose Alg. 3 mainly due to its success in training deep networks as observed in some prior works. Of course, there is room for improvement regarding this prediction process and can be the object of further research papers.

Numerical evaluation: For MNIST-image the test accuracy has been flat, and for the other two it is almost flat though slightly increasing. Yet, for illustrative purposes, the main idea is to show how faster our method is in the first epochs. The purpose of this method is not to achieve better generalization (i.e. reach better accuracy at convergence) but rather to show accelerated convergence compared to baselines. The learning rates have been tuned over a grid search and the best results have been reported. The choice of a constant learning rate was made to stick to our theoretical results. Runs with exponential decay or step decay can also be done for completeness.

Reviewer 4: We thank the reviewer for valuable comments and typos. Our response is as follows:

Numerical experiments: We focused on running different methods with same initialization for a fair comparison, and OPT-AMSGrad consistently outperforms vanilla AMSGrad with same initialization. As the reviewer kindly suggested, we will report error bars in the rebuttal version of the paper to show the advantage of OPT-AMSGrad with random runs in practice. We agree with the fact that our method empirically shows on Figure 2 and 3 better training performances (both in terms of loss and accuracy) but we must note how comparable and most of the time better than the baselines our method behaves at testing phase. • Typos and unclear figures will be changed according to the additional feedbacks.

• Figure 1 ought to represent the two-stage method and will be simplified in a linear manner as suggested. • AMSGrad and OPT-AMSGrad can indeed be written jointly (to highlight our contribution we decided to separate them). • The termination number T is classical to derive theoretical results in nonconvex optimization (see [Ghadimi & Lan, 2013]). • The bracket notation in **H3** is the inner product notation, this will be specified.