

1 We sincerely thank the four reviewers for their valuable feedback.

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

3 **Comparison with OPT-FTRL:** As stated in our introduction, we believe that the optimistic estimates presents an
4 advantage over vanilla AMSGrad like OPT-FTRL presents over vanilla FTRL. Keeping in mind that the overall goal of
5 our paper is to introduce a new method to solve a stochastic optimization problem where the objective function is a
6 (large) finite-sum, OPT-FTRL is not an option, hence our motivation for developing a counterpart of OPT-FTRL.

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

8 **Novelty of the contribution:** Although combining gradient prediction to AMSGrad update seems natural, as pointed
9 out in the first paragraph of Section 3, we would like to stress on how the embedding of the prediction process
10 (represented Figure 1) led to the two-stage algorithm OPT-AMSGrad (unlike the sequential structure of the original
11 AMSGrad) where, first an auxiliary variable \tilde{w} is updated and then the global model w . Also, as discussed in the
12 paper, optimistic learning is typically used in two-player-games, which is an online learning problem, and to the best
13 of our knowledge, this is the first proposal to apply optimistic acceleration to stochastic optimization problems (e.g.
14 training deep neural networks). Introducing the online optimization framework is a natural way to introduce our method,
15 providing related literature on optimistic methods.

16 **Comparison with other gradient prediction methods:** Comparing the way we predict the gradients is indeed
17 important in our study. We have devoted in Section F of the appendix an illustrative example of how this process can
18 impact the performances of our method.

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

20 **Gradient prediction algorithm:** We will add more explanations on why the average of the last gradients can be a
21 good approximation of the next one. While this may be counterintuitive, we invite the reviewer to read the citation we
22 make regarding our extrapolation method: "Regularized nonlinear acceleration" by Scieur, d'Aspremont and Bach,
23 NIPS 2016. We chose the latter reference mainly due to its success in training deep networks as observed in some prior
24 works. Of course, there is room for improvement regarding this prediction process for future research projects.

25 **Numerical Experiments:** The learning rates have been tuned over a grid search for all methods and the best per-
26 formances over 5 repetitions have been reported. The main motivation behind those plots is to show that adding an
27 optimistic update to the AMSGrad actually speed up the convergence in terms of both losses and accuracies. Given the
28 well-known advantages of Adam-type methods such as ADAM or AMSGrad, we did not compare to slower methods.

29 **Reviewer 5:** We thank the reviewer for valuable comments and typos. Our response is as follows:

30 **Discussion on the bounds:** **TO COMPLETE**

31 **Global convergence analysis:** The term *Global* is employed in the sense that it does not restrict the initialization of the
32 algorithm and our bound is true for any iteration (finite-time). In other words the result is global since it is true for any
33 initial point. Of course this is not related to the stationary point, as the objective function is nonconvex, no guarantees
34 are given regarding the nature of the obtained stationary point.

35 **Numerical experiments:** There are several works considering applying Alg. 3 in deep learning, e.g. [Nonlinear
36 Acceleration of Deep Neural Networks, Scieur et al., 2018], with positive results. As noted in their paper, in practice
37 extrapolation on CPU is faster than a forward pass on mini-batch and can be further accelerated on GPU. Moreover, note
38 that at each iteration, we only change one past gradient, so we do not need to compute the whole linear system every
39 time leading to practical efficiency. Secondly, the main focus of our paper is essentially the framework of integrating
40 optimistic learning with AMSGrad. We chose Algorithm 3 mainly because of the empirical success reported in prior
41 works. The choice of gradient prediction method is actually flexible. So, OPT-AMSGrad will definitely benefit from an
42 algorithm with faster running time and good prediction quality.

43 **Reviewer 6:** We thank the reviewer for valuable comments and typos. Our response is as follows:

44 **OPT-ADAM:** OPT-Adam has been developed in the particular case of an online problem, where the observations are
45 being presented in a streaming fashion. Here, our goal is to develop a method for the finite-sum stochastic optimization
46 problem. To the best of our knowledge, OPT-Adam has not been studied and no guarantees are given to claim that it
47 also performs well under this latter settings. Empirical evaluations actually show that OPT-AMSGrad is better.

48 **Boundedness of the iterates:** Lemma 2 is needed to ensure that H1 is verified. We will change the notations to avoid
49 any confusion. Lemma 2 indeed bounds the iterates of Algorithm 2 (OPT-AMSGrad) and should read $\|w_t^{(\ell)}\| \leq A_{(\ell)}$
50 where w_t are the weight estimates at iteration t . Then Lemma 2 holds for all iteration $t > 0$.