

1 We appreciate the valuable comments and positive feedback from the reviewers. We will carefully revise the paper  
2 accordingly to incorporate the comments.

3 **Reviewer #1: (Stepsize and preset  $T$ .)** Following the current analysis, for a general stepsize  $\eta_t$ , the convergence of  
4 stochastic update requires  $(\sum_{t=1}^T \eta_t^2)/(T \cdot \min_{t \leq T} \eta_t) \rightarrow 0$  and  $T \cdot \min_{t \leq T} \eta_t \rightarrow \infty$  as  $T \rightarrow \infty$  to handle the  
5 variance of stochastic semigradient. Thus, a diminishing adaptive stepsize such as  $\eta_t = 1/\sqrt{t}$  would also work, but the  
6 convergence rate would then become  $O(\log T/\sqrt{T})$ , which is slightly slower than the  $O(1/\sqrt{T})$  rate in our paper. For  
7 the same reason, an absolute constant stepsize does not guarantee convergence, since it fails to satisfy the first  
8 requirement. In view of the two requirements, we use the stepsize  $\eta_t = 1/\sqrt{T}$  to obtain the fastest rate  $O(1/\sqrt{T})$ . We  
9 will add a corresponding discussion in the revision.

10 **(Average of iterates.)** In the current analysis, the convergence rate is implied by the upper bound of a telescope sum  
11 (line 618 of the full paper). Without averaging the iterates, no convergence rate is available. Although the iterates using  
12 population semigradient would still converge, stochastic semigradient might cause divergence. Such a situation is  
13 analogous to convex optimization without strong convexity, where averaging the iterates is necessary [1].

14 **Reviewer #2: (Two-layer neural network.)** In this paper we consider neural network with one hidden layer. It is  
15 called a two-layer neural network following the recent line of work (e.g., [2]), since there is also a linear output layer.  
16 We recognize the potential confusion in terminology and will explicitly clarify that we mean a two-layer net with a  
17 single hidden layer.

18 **(Motivation for choosing the architecture.)** Such a shallow structure helps to characterize the learning dynamics and  
19 illustrate the connection to linear model with random features. With one hidden layer, it is already quite challenging to  
20 analyze the effects of using overparametrized neural networks for function approximation in RL.

21 **(Generalization to more complex networks.)** The results can be readily generalized to deep neural networks  
22 (multiple hidden layers with width  $m$ ) given the activation function is sufficiently smooth (e.g., sigmoid activation) and  
23 each layer is coupled with a suitable scaling factor. However, the ReLU activation used in this paper does not directly  
24 satisfy the smoothness requirement and therefore requires more delicate analysis.

25 **(MSPBE with oblique projections.)** Thanks for bringing up the oblique projection view. We will add a corresponding  
26 discussion in the revision. In the oblique projection paper, the difference between temporal difference-based and  
27 Bellman residual-based approaches arises due to the limited representation power of finite-dimensional linear function  
28 approximation. In comparison, overparametrized neural networks represent a larger infinite-dimensional function class,  
29 which alleviates the issues caused by limited representation power and therefore bridges the gap between the two  
30 approaches. In particular, Proposition 4.7 shows that neural TD attains the global minimum of MSBE (without the  
31 projection in MSPBE) under slightly stronger conditions.

32 **(State assumption.)** Our proof only relies on the fact that  $x$  is bounded, while the unit-norm assumption is used to  
33 simplify the presentation. An alternative view of this assumption is that the neural network has an additional (fixed)  
34 input layer that projects or embeds the “raw input”  $(s, a) \in \mathcal{S} \times \mathcal{A}$  to the unit sphere.

35 **(Reward assumption.)** Thanks for pointing this out. Coercive reward indeed requires more delicate analysis and is  
36 beyond the scope of this paper. We will revise the “without loss of generality” claim in the revision.

37 **(Function class  $\mathcal{F}_{B, \infty} - \hat{Q}(\cdot; W(0))$ .)** For any function class  $\mathcal{F}$  and function  $f'$ , the function class  $\mathcal{F} - f'$  is defined  
38 as  $\{g = f - f' : f \in \mathcal{F}\}$ . We will clarify this notation in the revision.

39 **(Minor comments.)** Thanks for pointing out the issues on notation and clarity. We will fix them in the revision.

40 **Reviewer #3: (One-point monotonicity.)** See line 591 of the full paper (deferred due to space limit) for more details  
41 on the notion of one-point monotonicity. We will move this to the main text in the revision as it is an important concept  
42 for this paper. Thank you for pointing this out.

43 **(Constants  $c_1, c_3, \nu$ .)** The exact polynomial dependency on  $c_1$  and  $c_3$  in the convergence rate is quantified in Lemma  
44 A.2 and the proof of Lemma E.2 (line 802) of the full paper (deferred due to space limit), which is omitted in big- $O$ 's  
45 when the lemmas are invoked. Meanwhile, the dependency on  $\nu$  is quantified in the proof of Theorem 5.3 (inequality  
46 (E.23) of the full paper) and is of order  $O(1/\nu)$ . We will move the dependencies to the main text in the next version.

47 **(How width affects rate.)** The effect of overparametrization is explicitly quantified in Theorems 4.4, 4.6, and 5.3 by  
48 the terms that decay with  $m$ , which denotes the width of the neural network. Roughly speaking, the convergence rate  
49 takes the form of  $1/\sqrt{T} + 1/\sqrt{m}$ . As  $m \rightarrow \infty$  (or at least  $m = \Omega(T)$ ), the rate reduces to  $1/\sqrt{T}$ , where  $T$  is the  
50 number of iterations. In other words, the “error of implicit linearization” diminishes as the neural network has more  
51 parameters. We will include a discussion of how width affects the convergence rate in the next revision.

52 [1] Bubeck, S. (2015). Convex optimization: Algorithms and complexity. Foundations and Trends in Machine  
53 Learning, 8 231–357.

54 [2] Arora, S., Du, S. S., Hu, W., Li, Z. and Wang, R. (2019). Fine-grained analysis of optimization and generalization  
55 for overparameterized two-layer neural networks. arXiv preprint arXiv:1901.08584.