

1 We appreciate the valuable comments from the reviewers. We will revise accordingly.

2 **Reviewer #1. (Novelty in techniques.)** Our proof of the convergence of PPO/TRPO has several building blocks,
3 which, to the best of our knowledge, is not covered by existing work.

4 • **TD and SGD:** Our unified analysis covers both TD and SGD, where **TD is not covered by any existing supervised**
5 **learning analysis**. In particular, the semigradient used in the TD update is **not an unbiased stochastic gradient of**
6 **any supervised learning objective**. To see this, one may take the derivative of the semigradient, which gives an
7 asymmetric matrix. In contrast, in supervised learning, such a matrix is the Hessian, which must be symmetric.

8 A “straightforward adaptation” of existing supervised learning analysis does not yield the global convergence of TD. In
9 fact, most algorithms in supervised learning are known to at least converge (although they may converge to undesired
10 stationary points). In contrast, TD with nonlinear function approximator **is known to generally diverge [1]**, which
11 eludes existing supervised learning analysis. Our unified analysis of both TD and SGD shows how overparametrization
12 allows bypassing the divergence of TD, which has not been observed in the context of supervised learning before.

13 • **Nonconvex mirror descent:** Most existing analysis of mirror descent’s convergence to a **global optimum** builds
14 on the critical assumption that the objective is convex, which is not the case for the objective $J(\pi)$ in RL. Our **global**
15 **convergence** analysis of PPO/TRPO only builds on **one-point monotonicity** (Lemma 5.1), instead of convexity. In
16 fact, we are not aware of any existing analysis of **infinite-dimensional** mirror descent under one-point monotonicity.

17 • **Error propagation:** RL is divided into **policy-based** and **value-based** approaches. **PPO/TRPO falls into the**
18 **former**, while **Q-learning falls into the latter**. More specifically, PPO/TRPO explicitly tracks a policy, while Q-
19 learning does not. In particular, the Q-function tracked in Q-learning is not the action-value function of any policy. As
20 a result, the error propagation in Q-learning relies on the **γ -contraction of the Bellman optimality operator**, which
21 operates on the **Q-function**. In contrast, the error propagation in PPO/TRPO in our analysis relies on the **convergence**
22 **of mirror descent**, which updates the **policy**. Such two mechanisms are not comparable.

23 **(Expectation over initialization.)** Yes, the expectation is also taken over the initialization of the neural networks. Our
24 notation indeed lacks clarification. Thanks for pointing this out. We will revise accordingly.

25 **(Sufficiently large function class.)** Yes, the divergence is largely due to the limited representation power of **finite-**
26 **dimensional** linear function approximators. In contrast, overparametrized neural networks provide sufficient represen-
27 tation power as **infinite-dimensional** nonlinear function approximators, which enables the global convergence.

28 **Reviewer #2. (Shallow architecture.)** Thanks for pointing this out. The name “two-layer neural network” is standard
29 in recent work (see, e.g., [2]), counting the output layer as one of the layers. We will emphasize more on the shallow
30 architecture considered here.

31 **(Bounded reward.)** The state-value function and action-value function defined in Section 2 are normalized by a factor
32 of $1 - \gamma$. Hence, they have the same upper bound R_{\max} as the reward function.

33 **(Simulation.)** Thanks for the suggestion. We will add illustrative simulation examples to the appendix.

34 **Reviewer #3. (Conclusion section.)** We will add a conclusion section to briefly summarize this paper as well as
35 potential future work.

36 **(Nonconvexity, infinite-dimensionality, and technical challenges.)** The nonconvexity arises from two aspects: (i)
37 The neural network parametrization of the energy-based policy and the action-value function makes the subproblems of
38 policy improvement and policy evaluation nonconvex. (ii) The RL objective $J(\pi)$ is also nonconvex in π . Meanwhile,
39 the continuous state space introduces infinite-dimensionality. The nonconvexity and infinite-dimensionality combined
40 makes the global convergence analysis challenging.

41 **(Algorithm 1.)** Yes, in Algorithm 1, lines 5 and 6 are just one projected TD step. Correspondingly, lines 5 and 6 in
42 Algorithm 2 are one projected SGD step.

43 **(Overparametrization.)** Our analysis does require the width of the neural network to be large enough. In practice, in
44 order to ensure the desired representation power of the neural network, a larger number of parameters is required.

45 [1] TSITSIKLIS, J. N. and VAN ROY, B. (1997). Analysis of temporal-difference learning with function approximation.
46 In Advances in Neural Information Processing Systems.

47 [2] MUNOS, R. and SZEPESVARI, C. (2008). Finite-time bounds for fitted value iteration. Journal of Machine
48 Learning Research.

49 [3] ARORA, S., DU, S. S., HU, W., LI, Z. and WANG, R. (2019). Fine-grained analysis of optimization and
50 generalization for overparameterized two-layer neural networks. arXiv:1901.08584.