

1 **Reviewer 1: Q:** Compare and highlight new challenges of analyzing TDC relative to other GTD algorithms in [30,34].

2 **A:** (a) As mentioned in [14], TDC does not have an explicit saddle point representation as GTD and GTD2, and hence  
3 its analysis cannot follow the convex-concave optimization framework developed in [30,34]. (b) [30,34] assume that  
4 two variables' updates have the same constant stepsize. For TDC, we analyze more general cases: two stepsizes have  
5 different diminishing rates, and two stepsizes are different valued constants. Consequently, in our case, interaction  
6 between two variables requires more sophisticated techniques to analyze, e.g., recursively sharpening error bounds.

7 **Q:** The paper generalizes stagewise stepsize in conventional (one timescale) optimization to two-timescale optimization.  
8 Discuss new challenges of analyzing algorithms with blockwise diminishing stepsize in this new settings.

9 **A:** Comparing to conventional optimization with stagewise stepsize, here we need to handle the bias induced by non-i.i.d.  
10 samples and characterize non-asymptotic behavior of two timescale variable update. Hence, the update scheme for  
11 stepsize and block length in each block is designed based on two time-scale analysis to yield linear convergence rate  
12 blockwisely and desirable sample complexity.

13 **Q:** How the theoretical guarantees can be affected in the non-asymptotic analysis of actor-critic and gradient Q-learning.

14 **A:** Since both actor-critic and gradient Q-learning algorithms are two time-scale algorithms, our non-asymptotic analysis  
15 for two time-scale algorithms can be very useful. Moreover, analysis of these two algorithms will further require to deal  
16 with their special structures such as policy update, presence of multiple fixed points, local convergence, etc.

17 **Reviewer 2: Q:** How is  $\theta^*$  obtained in the experiments.

18 **A:** In our experiment (Garnet problem), since we pick behavior policy  $\pi_b$  and transition probability  $p(s'|s, a)$ , the  
19 stationary distribution  $\mu_{\pi_b}$  can be computed. Since we also know target policy  $\pi$  and feature matrix  $\Phi$ , we can compute  
20 the matrix  $A$  and the vector  $b$  by definition to obtain  $\theta^* = -A^{-1}b$ . Alternatively,  $\theta^*$  can also be estimated by running  
21 the algorithm with diminishing stepsize for sufficiently long time and taking the average of outputs of several runs.

22 **Q:** How would worst-case errors predicted by the bound compare to errors observed empirically in experiments.

23 **A:** Our theory captures how the error bound changes with stepsize diminishing parameters  $(\nu, \sigma)$ , which agrees with  
24 how the empirical error changes with stepsize diminishing parameters in our experiments (see Fig. 1 in the paper).  
25 Furthermore, specializing Theorem 1 to i.i.d. scenarios, our convergence rate order-wisely matches the best known  
26 result in [Dalal et al. COLT 2018]. It is a good idea to plot theoretical errors and compare with empirical bounds. The  
27 main challenge here is that precisely estimating some parameters in the error bound (e.g., eq (25)) can be difficult  
28 (although they are known to be constants in convergence analysis). For example, mixing time parameters  $\tau_\alpha$  and  $\tau_\beta$  in  
29 (25) depend on geometric ergodicity of Markov chain, but constants  $m$  and  $\rho$  (see Assumption 3) are usually difficult to  
30 estimate in practice. We are currently further exploring such an issue.

31 **Q:** Besides implications for the choice of step-size, do these bounds provide insight on what properties of the problem,  
32 the behavior policy, and the representation affect the rate of convergence?

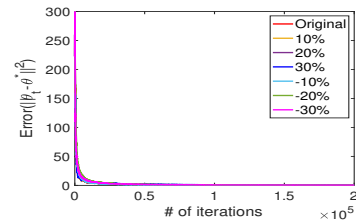
33 **A:** Theorem 1 (more precisely eq (25) in suppl.) captures how other properties (besides stepsize) affect convergence  
34 rate. For example, convergence rate depends on  $\lambda_\theta$ , which is lower-bounded by the largest eigenvalue of matrix  
35  $2A^T C^{-1} A$ , and such a matrix is determined by behavior policy  $\pi_b$ , target policy  $\pi$ , transition probability  $p(s'|s, a)$  and  
36 feature matrix  $\Phi$ . Convergence rate also depends on the mixing time  $\tau_\alpha$  due to geometric ergodicity of the Markov  
37 chain, which is determined by  $\pi_b$  and  $p(s'|s, a)$ . Other constant terms in (25) such as  $L_{f_1, \theta}$ ,  $K_{f_1}$  and  $K_{g_1}$  capture the  
38 dependence on  $\pi_b$ ,  $\pi$ ,  $\Phi$  and the discount factor  $\gamma$ .

39 **Q:** Explain what "more flexible" mean when saying gradient TD are "more flexible than on-policy learning in practice."

40 **A:** We meant gradient TD algorithms are flexible because they converge even with off-policy data and hence can exploit  
41 abundant samples (obtained under behavior policies) for learning when the on-policy samples are limited.

42 **Reviewer 3: Q:** How to set blocksize properly without prior knowledge and how  
43 robust the algorithm is with respect to blocksize hyperparameter.

44 **A:** In practice, blocksize  $T_s$  and stepsize  $\alpha_s$  are set by parameter tuning, but  
45 we do not directly tune them for all blocks because there are too many tuning  
46 parameters this way. Instead, Theorem 3 indicates that  $T_s$  and  $\alpha_s$  for all blocks  
47 are fully determined by only four parameters  $\epsilon_0$ ,  $|\lambda_x|$ ,  $C_7$ , and  $\eta$ , and among  
48 them  $|\lambda_x|$  and  $\eta$  can be estimated by matrices  $A$  and  $C$  from samples. Hence  
49 we mainly tune only  $\epsilon_0$  and  $C_7$ . Our experiments demonstrate that this approach yields desirable performance. For  
50 robustness, we run experiments (see the figure on the right) and find that perturbing blocksize even by  $\pm 30\%$  for all  
51 blocks changes the convergence rate only very slightly, demonstrating that the performance of algorithm is very robust  
52 to blocksize.



50 robustness, we run experiments (see the figure on the right) and find that perturbing blocksize even by  $\pm 30\%$  for all  
51 blocks changes the convergence rate only very slightly, demonstrating that the performance of algorithm is very robust  
52 to blocksize.