

1 **Reviewer #1: Q1.** *The projection step is impractical. In addition, if one notices, later in the proof of Lemma 2 (see eq.*  
2 *(27) in the Appendix), another assumption requires the norm of feature vectors to be less than 1. Making these two*  
3 *assumptions simultaneously, does not necessarily guarantee that the true Q-function can be approximated by  $\phi^T \theta$ .*

4 **Response:** As we clarify in our answer to Q2 below, projection is not needed for practically implementing SARSA,  
5 and convergence still holds. Even if it is used, to guarantee approximation of Q-function by  $\phi^T \theta^*$ , as stated in our  
6 Theorem 1, we require that the optimal  $\theta^*$  is within the projection radius. Such a projection radius can be obtained in  
7 practice by using our Lemma 1 and designing an online estimator of  $w_l$ .

8 **Q2.** *Later, it is claimed that without the projection, the method is still useful, and one just needs to set R to possibly a*  
9 *very large value; however, it can be seen from the right-hand side of equations (3), (4), and (11), that the second and*  
10 *third power of R shows up through G, and thus the bounds will possibly be very loose.*

11 **Response:** We clarify that projection is NOT necessary to implement SARSA and does not affect its convergence based  
12 on existing literature. (Our statement "set R to possibly a very large value" was not accurate.) [Gorden 2001] has  
13 shown that SARSA converges to a bounded region. [Melo et. al. 2008] and [Perkins and Precup 2003] further showed  
14 that SARSA with a Lipschitz continuous policy improvement operator converges even without projection. Thus  $\theta_t$   
15 generated by SARSA is already bounded by itself, and projection is not necessary. Thus, the non-asymptotic bound for  
16 SARSA without projection is the same as if we use projection with  $R = \max_t \|\theta_t\|_2$  (only for the purpose of analysis).  
17 Moreover, R is determined by the nature of algorithm. If it happens to be large, then the error should be large by nature,  
18 which doesn't mean the bound is loose.

19 **Q3.** *Another concern was also raised by [Chen et. al. 2019]. It is claimed in [Chen et. al. 2019] that Thm 1 of [Melo et.*  
20 *al. 2008], that this work builds on, cannot be verified. This concern further weakens the contributions of this work.*

21 **Response:** We clarify that this paper studies SARSA, **not** Q-learning. Our proof for SARSA does not have the issue  
22 pointed out in [Chen et. al. 2019] for Q-learning. We are aware of the fact that one step in the proof of Thm 1 in [Melo  
23 et. al. 2008] cannot be verified for Q-learning, which might be due to greedy policy and off-policy training taken in  
24 Q-learning. Such an issue does not affect SARSA as an on-policy algorithm.

25 **Q4.** *Assumption 2 requires C small enough to guarantee a negative eigenvalue. No clear characterization of  $w_s$ . In*  
26 *practice, how one can guarantee a "feasible" or "meaningful" policy for which it is possible to have such a small C.*

27 **Response:** In practice, one can numerically tune the parameter to find a "feasible" policy improvement operator, e.g.,  $\beta$   
28 in softmax (which corresponds to a C). Assumption 2 (for guaranteeing convergence) can be empirically checked to  
29 provide guidance for parameter tuning.

30 **Q5.** *It is also recommended to provide some numerical tests to verify the theoretical results.*

31 **Response:** We have run numerical results and will include them in the revision.

32 **Reviewer #2: Q6.** *Write the proof for Lemma 5 explicitly.* **Response:** Done.

33 **Q7.** *Remind that Eq (50) follows from (49) because of the relationship between  $\alpha_t$  and  $w_s$ .* **Response:** Done.

34 **Q8.** *How realistic is assumption that for any  $\theta$  the induced Markov chain is ergodic (right before Assumption 1)? Is it*  
35 *guaranteed for any special class of environments or representations? Discussion on the plausibility of assumptions.*

36 **Response:** First note that a Markov chain is uniformly ergodic if it is irreducible (i.e., possibly get to any state from  
37 any state) and aperiodic [Levin & Peres 2017]. Now consider an environment for which there exists a policy that  
38 can map any state to any state with nonzero probability (i.e., irreducibility holds) and can get back to the same state  
39 aperiodically. (Note that such environments are commonly encountered in practice.) Then for any  $\theta$ , as long as the policy  
40 improvement operator explores (i.e., with non-zero probability to take any action at any state), the induced Markov  
41 chain remains to be irreducible and aperiodic, and is hence ergodic. Therefore, such an assumption is realistic, and is  
42 guaranteed for aforementioned environments. Please see response to Q1 and Q4 for other assumptions.

43 **Q9.** *Does this result bear implications for framing a RL problem or designing features? Highlight any insights that*  
44 *follow from these results.*

45 **Response:** First, our result characterizes sample complexity of SARSA for both constant and diminishing step sizes,  
46 which is useful for choosing learning rate to design fast RL algorithms. Second, our result indicates that the faster the  
47 underlying Markov process mixes, the faster SARSA converges. This motivates the design of tricks (e.g., experience  
48 replay for TD and Q-learning [Wang et. al. arXiv:1809.08926 2017]) to improve the mixing property of Markov process.

49 **Reviewer #3: Q10.** *The analysis is limited in linear function approximation case.*

50 **Response:** One great advantage of linear function approximation is that it is very easy to implement. Even for this  
51 case, there is no existing finite sample analysis in the literature, and this paper accomplishes this step (which is already  
52 technically nontrivial). It is of great interest to further explore more advanced function spaces, such as deep neural  
53 networks. Looking forward, this work serves as a first step towards understanding the more complicated case.