

1 **To R1, R2, R4: *About the proof sketch:*** Due to the 8 pages limit, we had to put the proof sketch in Appendix A (Page
2 12-14) and use the main text to highlight the contribution. In the camera ready version, we can reorganize and move the
3 proof sketch into the main text, and also add references/pointers to the specific section in the appendix.

4 **To R1: *The paper is hard to follow...*** We apologize we did not provide enough background. We will provide more
5 explanation to those terminologies in the revision.

6 ***Contribution compared with previous works:*** The main contributions of our work include: (1) We are the first to provide
7 a finite-time analysis for a practical actor-critic algorithm introduced in [25]; (2) Our analysis does not require the
8 unrealistic i.i.d. assumption and directly deals with the Markov decision process; (3) Our analysis provides a better
9 sample complexity $\mathcal{O}(\epsilon^{-2.5})$ than the best-known result $\mathcal{O}(\epsilon^{-4})$ in previous work under strong assumptions [16,21].

10 ***1. Why decoupled actor-critic assumption?*** The decoupled actor-critic is not an assumption, but a different algorithm
11 appearing in [16] and [21]. It is not a practical algorithm, but easier to analyze. Our two time-scale actor-critic algorithm
12 is more realistic and more sample efficient.

13 ***2. What is τ_t ?*** τ_t is the mixing time of the Markov chain, which characterizes the time it takes the ergodic Markov chain
14 in Assumption 4.2 to converge to its stationary distribution.

15 ***3. What is Markovian noise?*** Both actor and critic are updated based on observation tuples $\{O_t = (s_t, a_t, s_{t+1})\}_{t=0,1,\dots}$.
16 In previous work [16, 21], they assume that each tuple is sampled i.i.d. in order to simplify the analysis. However, this
17 is obviously not true in practice. In this paper, we follow [3] and directly deal with the true data which are sampled
18 from the Markov decision process. We refer to this setting as the ‘‘Markovian noise’’ setting.

19 ***4. How is ‘‘iterative refinement’’ used?*** The ‘‘iterative refinement’’ is not used in our paper. It is used in [36] and more
20 details can be found therein.

21 ***5. Proof of Corollary 4.9:*** The proof is in Section C.4, line 667.

22 ***6. Proof of Lemma B.3:*** This is a typo. All a_i should be a_k .

23 ***7. Proof of Lemma C.2*** The proof is as follows:

$$\begin{aligned} \|\Delta h(O, \eta, \omega, \theta)\|^2 &:= (\eta(\theta) - \eta + (\phi(s') - \phi(s))^\top (\omega - \omega^*))^2 \cdot \|\nabla \log \pi_\theta(a|s)\|^2 \\ &\leq [2(\eta(\theta) - \eta)^2 + 2((\phi(s') - \phi(s))^\top (\omega - \omega^*))^2] B^2 \\ &\leq [2(\eta(\theta) - \eta)^2 + 2\|\phi(s') - \phi(s)\|^2 \|\omega - \omega^*\|^2] B^2 \leq [2(\eta(\theta) - \eta)^2 + 2 \cdot 4 \cdot \|\omega - \omega^*\|^2] B^2, \end{aligned}$$

24 where the equality is by the definition of $\Delta h(O, \eta, \omega, \theta)$, the first inequality is by $(a+b)^2 \leq 2a^2 + 2b^2$ and Assumption
25 4.3(a), the second inequality is by Cauchy-Schwartz, and the last inequality is by triangle inequality and $\|\phi(s)\| \leq 1$.

26 **To R2: *Technical novelty in Remark 4.8:*** Take the estimation error \mathbf{z}_t as an example, the inequality at line 642 involves
27 bounding $\mathbb{E}\|\mathbf{z}_t\|^2$ with $\mathbb{E}\|\mathbf{z}_t\|$ at its right hand side. Directly unrolling the equation at line 642 yields a loose result and
28 a complicated proof, as done in [36]. We find it is viable to postpone the unrolling and compute the average estimation
29 error which can give a tighter bound. We will elaborate it in Remark 4.8 in the revision.

30 ***Technical difference with [Y]:*** Thanks for pointing out the related works which we were not aware of previously. We
31 will compare with them in the revision. The problem settings of both works are very different. In specific, [Y] considers
32 updates that are linear in the two parameters θ_t and ω_t . In contrast, the actor-critic updates in our paper is not a linear
33 function of θ_t and ω_t , i.e., the policy gradient update for θ_t . So the analysis in [Y] is not directly applicable to our
34 setting, and our analysis on the actor θ_t update requires a very different approach.

35 ***Where the first term in (4.3) and (4.4) come from?*** The first term in (4.3) comes from the term I_1 at line 648, which is
36 related to the estimation error of last iterate ω_t . Similarly, the term in (4.4) comes from I_1 at line 615, which is related
37 to the estimation error of the last iterate η_t for the average reward.

38 **To R3: *The ϵ -approximate stationary point:*** We will explicitly acknowledge the existence of function approximation
39 error to avoid any confusion.

40 ***Compatible function approximation:*** It is possible to use compatible function approximation instead of a fixed linear
41 function approximation. The potential difficulty is to efficiently estimate the Q function for a given state-action pair,
42 which might involve starting another sampling trajectory.

43 ***Discounted setting:*** It is possible to extend our analysis to discounted MDPs. As you suggested, we can discarding each
44 transition w.p. $1 - \gamma$ and restarting the episode. We will add a discussion in the future work section.

45 ***Q-function or Advantage function:*** Our analysis is applicable to both advantage function $\Delta(s, a)$ and Q-function, with a
46 very minor change in the analysis. We use the advantage function just following the convention of practice.

47 ***Regularized critic:*** Thank you for pointing out the related work [ICML2020] and suggesting this very promising idea.
48 We will comment on this work and study the regularized critic in our future work.

49 **To R4: *The proofs are extending over more than 20 pages and they are marked as ‘‘sketches’’:*** This is a misunderstanding.
50 To clarify, we actually have both sketches of proofs (Appendix A, pp. 12-14) and detailed proofs (Appendix C). So it is
51 not the sketch that is lengthy.

52 ***About Theorem 4.5*** The first term is the linear approximation error; the second term is from upper bounding the
53 performance function; the third term is due to the stochastic variance and Markovian noise; and the last term is the
54 critic’s error. More details can be found in the proof sketch, at line 446.

55 ***Joint loss:*** In Open AI’s implementation of A2C, the gradients of the joint loss w.r.t. the actor and the critic can actually
56 be separated, so our analysis still holds.