

1 **Reviewer #1** We thank the reviewer for the valuable comments.

2 1. Large constants in Ψ and κ :

3 The constant 400 in Ψ comes from Lemma B.5, which shows how to construct a normal barrier in the lifted
4 domain from an arbitrary self-concordant barrier in the original decision set. However, as we mention in
5 Footnote 2, our algorithm works with *any* normal barrier, not just this particular one (since our proofs only use
6 those general properties of normal barriers from Section B.2). Therefore, the constant could be much smaller
7 as long as a normal barrier with a smaller constant exists, which is indeed true for several canonical examples
8 (see Nesterov and Nemirovskii, 1994). Similarly, the constant in κ is determined by the constant in Ψ (see the
9 proof of Lemma B.12), and thus could be smaller as well.

10 As we discuss in the paper, our algorithm is essentially SCRiBLE with a new sampling scheme and a new
11 adaptive learning rate schedule. We thus believe that our algorithm is at least as empirically viable as SCRiBLE.

12 2. Other minor comments and missing related works:

13 Thanks for pointing these out! We will revise the paper accordingly.

14 **Reviewer #2** We thank the reviewer for the valuable comments.

15 **Reviewer #3** We thank the reviewer for the valuable comments.

16 1. Efficiency of OMD with log-barrier regularizer for MDPs:

17 In general, since the OMD step is a convex optimization problem with $\text{poly}(|S|, |A|)$ linear constraints, one
18 could apply any standard convex solver to implement the algorithm efficiently. Jin et al. 2020 were able to
19 reduce the problem to another optimization problem with only positivity constraints due to the special structure
20 of the entropy regularizer (which does not hold for log-barrier), but in the end they still require applying a
21 convex solver to implement the OMD step.

22 2. How to sample s_t efficiently:

23 Yes, the way you mention is correct and efficient (except that in the end one also needs to normalize the norm
24 to 1). Another efficient way to sample s_t is to first uniformly randomly sample a point s from a d -dimensional
25 unit sphere, then let $s_t = H_t^{\frac{1}{2}} \begin{bmatrix} s \\ 0 \end{bmatrix}$, and finally normalize s_t . We will add these discussions to the final version.

26 3. Typos and clarification:

27 Thanks for pointing out the typos on Page 13 and Page 16. We will revise the paper accordingly. As for
28 the question regarding line 587, note that in the description of Algorithm 2, we restrict the choice of w_t in
29 $\Omega' = \{w = (w, 1) : w \in \Omega, \pi_{w_1}(w) \leq 1 - \frac{1}{T}\}$, therefore, we have $1 - \pi_{w_1}(w_t) = 1 - \pi_{w_1}(w_t) \geq \frac{1}{T}$.

30 **Reviewer #4** We thank the reviewer for the valuable comments and suggestions, and will take them into account
31 when revising the paper.