1 We thank all reviewers for very helpful comments. This letter addresses several major questions raised by the reviewers.

2 **A common question: numerical evaluation.** We will add a series of nu-
3 merical experiments to demonstrate the minimax optimality of the model-
4 based approach studied herein. Fig. 1 depicts some numerical plots of this
5 kind. Here, we adopt a (least favorable) example designed in the minimax
6 bound in Azar et al. [2] (and also Wainwright [37]), which primarily fo-
7 cuses on the effect of the discount complexity $\frac{1}{1-\gamma}$ on the sample complex-
8 ity. More specifically, consider the MDP $(\mathcal{S}, \mathcal{A}, P, r, \gamma)$, where $\mathcal{S} = \{0, 1\}$,
9 $\mathcal{A}_0 = \{0\}, \mathcal{A}_1 = \{0, 1\}, P(0|0, 0) = 1, P(1|1, 0) = p, P(1|1, 1) = q$ and
10 $r(0, 0) = 0, r(1, 0) = r(1, 1) = 1$. According to Azar et al., the quantities $p$
11 and $q$ are taken to be $p = \gamma + \frac{2\gamma(1-\gamma)^2\varepsilon}{(1+\gamma)^2}$ and $q = \gamma - \frac{2\gamma(1-\gamma)^2\varepsilon}{(1+\gamma)^2}$. As illustrated
12 in Fig. 1, the numerical sample complexity per state-action pair $N$ scales
13 on the order of $\frac{1}{(1-\gamma)^3\varepsilon^2}$ for varying choices of accuracy level $\varepsilon$, which is
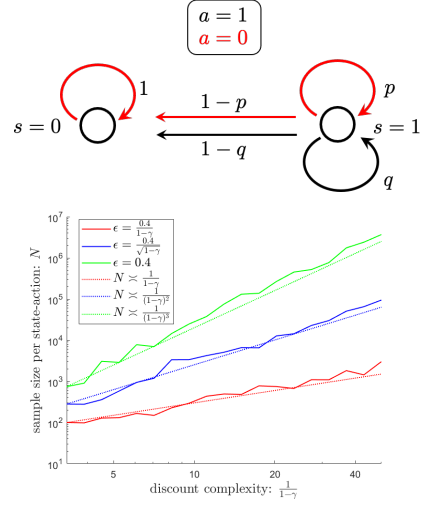14 consistent with our theory.



Figure 1: Numerical experiments.

15 **Specific questions by Reviewer 1:** 1. *Necessesity of reward perturbation.*
16 Indeed, reward perturbation is introduced merely to facilitate analysis. Given
17 that we recommend extremely small perturbations, the difference between the
18 perturbed and original MDPs is almost unnoticeable in the experiments. While
19 our current analysis does not work without reward perturbation, it would be
20 important to see whether this can be removed with more refined analysis.

21 **Specific questions by Reviewer 2:** 1. *Empirical illustration of minimax optimality.* We have conducted some
22 illustrative numerical experiments to address this comment; see the response above for "Numerical evaluation".

23 2. *Choice of $\alpha$.* Indeed, our theory and analysis hold for any constant $\alpha > 1$. The current paper picks the specific choice
24 $\alpha = 5$ only to convey that the perturbation $\frac{(1-\gamma)\varepsilon}{|\mathcal{S}|^\alpha|\mathcal{A}|^\alpha}$ can be very small; we shall make it more clear in the final version.

25 3. *Why the sample size barrier $N > 1/(1-\gamma)^2$ appeared in Agarwal et al. [1].* Take Section 4.3 of the Arxiv version
26 of Agarwal et al. for example: the contraction factor $\gamma\sqrt{\frac{8\log(|\mathcal{S}||\mathcal{A}|/(1-\gamma)\delta)}{N}}\frac{1}{1-\gamma}$ needs to be smaller than 1, which
27 requires the sample size per state-action pair to exceed $N > 1/(1-\gamma)^2$ (up to log factor). We shall explain it in the
28 final paper.

29 **Specific questions by Reviewer 3:** 1. *Motivation for perturbed rewards and paper structure*: Thanks for the helpful
30 suggestion. We will elucidate the motivation and intuition of reward perturbation earlier on in the revised paper.

31 2. *Empirical evaluation*: We have conducted some illustrative numerical experiments which will be added to the revised
32 paper; see the response above on "Numerical evaluation".

33 3. *Paper structure and broader impacts*: We will restructure the appendix (so as to reduce repetition) and add new
34 statements about potential societal impacts as suggested by the reviewer.

35 **Specific questions by Reviewer 4:** 1. *Constraint on accuracy level $\varepsilon$.* In fact, our result (see Theorem 1) holds for an
36 arbitrary choice of $\varepsilon \in (0, \frac{1}{1-\gamma}]$. This accommodates an arbitrary $\varepsilon$, and there is no lower bound on the accuracy level $\varepsilon$
37 in our main theorem. To be a bit more specific, our sample complexity bound reads $\frac{|\mathcal{S}||\mathcal{A}|}{(1-\gamma)^3\varepsilon^2}$, while prior work like
38 Agarwal et al. reads $\max\left\{\frac{|\mathcal{S}||\mathcal{A}|}{(1-\gamma)^3\varepsilon^2}, \frac{|\mathcal{S}||\mathcal{A}|}{(1-\gamma)^2}\right\}$ (up to some log factors). Here, the red text represents the sample size
39 barrier, and is removed in our theory. We understand from the reviewer's comment that there might be confusion in our
40 current discussions/remarks, and we shall rephrase them in the final paper to clarify any such confusion.

41 2. *The role of reward perturbation and choice of $\alpha$.* Reward perturbation is introduced mainly to break the ties (so as
42 to ensure uniqueness of optimal policy). Fortunately, even an extremely small level of perturbation suffices for this
43 purpose (as long as it is not exponentially small). Our analysis and theorems hold for any fixed constant $\alpha > 1$. Here,
44 we pick $\alpha = 5$ only to emphasize that a fairly small level of reward perturbation suffices for our analysis to work (more
45 specifically, a fairly small level of perturbation suffices in breaking the ties). This will be made clear in the final paper.

46 3. *About perturbed reward samples.* Here, reward perturbation is only enforced when running the planning algorithms
47 on empirical MDPs (basically, we collect the true reward samples $r$, but use $r_{\rm p} = r + \zeta$ in the algorithm). In other
48 words, it is not necessary to collect new samples for this purpose; we shall clarify this point in the revision.

49 4. *Missing state-action dependency in Eq. (10).* Here, $N$ is defined as the sample size *per state-action pair* (so that the
50 total sample size should be $|\mathcal{S}||\mathcal{A}|N$. We shall make it more clear in the revision.