

1 Adaptive Experimental Design with Temporal Interference: A Maximum Likelihood Approach

2 Response to reviewers

3 We are grateful to the referees for their thoughtful comments regarding our paper. Regarding typos and other detailed
4 suggestions, we plan to incorporate those prior to submission of our camera ready version. Below we have also provided
5 feedback in response to major comments and requests from the referees.

6 **Practical considerations.** Several referees commented on practical considerations (Reviewers 1, 3, and 4): in particular,
7 the performance of maximum likelihood estimation (MLE) with large state spaces and on finite horizons, as well as its
8 computational complexity. More broadly, we note that in followup work, we have developed an alternative approach to
9 experimental design for temporal interference using sample average estimation (SAE). In this work, we study sampling
10 strategies based on the *regenerative method*; together, SAE and regenerative policies ensure consistency while being
11 practically implementable. The idea is that we commit to a single state (the *regeneration state*), and only allow ourselves
12 to switch chains in this state. This method can be used together with sample average estimation (SAE) of rewards,
13 rather than (MLE). Further, it requires no advance knowledge of the state space, other than the regeneration state. This
14 is a much more practical, scalable solution; although it is not as sample efficient as the policy and MLE in our present
15 paper, we can use similar techniques to obtain the optimal regenerative policy with SAE. Though we could not include
16 this work due to space constraints, we will add some discussion of this extension to the paper to address concerns
17 regarding practical implementation.

18 **Reviewer 1.** Regarding how to interpret our results, note that consistency follows under very general conditions: *any*
19 sampling procedure that samples both experiments infinitely often in each state will ultimately yield consistent estimates
20 via the MLE. By using the MLE instead of just averaging rewards obtained in runs of each chain, we avoid temporal
21 interference completely. The main contribution of our paper is to provide a strong characterization of the *sample*
22 *efficient* experimental design for the MLE.

23 Regarding practical considerations and state space complexity, note that our theory shows that for any finite state space
24 our policy eventually outperforms any other TAR policy. Nevertheless, you are right that for a given sample size, our
25 policy may be computationally complex; see our comments above on practical considerations.

26 Regarding multiple treatments, we conjecture that the optimal policy in the multiple treatments setting eventually looks
27 essentially like ours, once the best and second best treatments have been identified.

28 **Reviewer 2.** Thank you for your review. Regarding making the paper more accessible to the Neurips community: we
29 will plan to provide some more intuition for the main results in the final version.

30 **Reviewer 3.** Regarding the convex optimization problem in Theorem 13 and Section 6, we expect that in many
31 applications the switching time is slow enough relative to computation time, so that standard convex optimization
32 techniques can be employed. That said, we agree that computational complexity is an important practical issue; see our
33 discussion above.

34 The set \mathcal{K} is a closed, compact, convex polytope in $2|S|$ -dimensional space (where S is the state space).

35 Regarding what policies are time-average regular (TAR), we emphasize that TAR is a weak regularity requirement: it
36 says that the fraction of time steps in which chain ℓ is sampled in state x converges to a well-defined random variable
37 (possibly deterministic). Virtually any reasonable policy will satisfy this requirement. Policies which, e.g., switch
38 chains on exponentially increasing timescales will not be TAR, but such policies are not likely to be used in practice.

39 **Reviewer 4.** Regarding more general optimality results, we agree with your conjecture; this remains an important open
40 direction. Note that for any fixed TAR policy, the MLE is asymptotically efficient given the samples collected by that
41 policy. Further, our results show how to compute the optimal TAR policy when estimating via the MLE.

42 Regarding obtaining higher reward by combining the two chains, our work is entirely focused only on *estimation* of the
43 difference in steady state rewards of the two chains, rather than *optimization* of the cumulative reward obtained. An
44 interesting question concerns whether regret optimal policies can be designed to maximize the cumulative reward using
45 only the two chains; this is an interesting reinforcement learning problem for future study.

46 **Other comments: Synthetic evaluation.** We agree with Reviewers 1 and 3 that synthetic evaluation would be valuable
47 to add to our paper. Synthetic evaluation would primarily be valuable to study finite horizon performance, as our paper
48 provides a full characterization of optimal asymptotic sample efficiency. We also emphasize that our paper is primarily
49 a theoretical study of optimal experimental design in this setting; practical considerations can lead to different preferred
50 designs (cf. our discussion above). Intuitively we believe the regenerative method leads to designs with improved finite
51 horizon performance; we plan to carry out a synthetic study to compare and contrast finite horizon performance of
52 various designs as part of our future work.