

1 We thank the reviewers for their thoughtful feedback. We believe that a few misreadings of our work made some of the
2 evaluations overly harsh and would ask reviewers to reconsider our paper in the light of clarifications provided below.

3 **1. Why Online learning? (R2)** It is acknowledged in the literature (Chakraborty and Murphy, 2014; Chakraborty
4 and Moodie, 2013) that most of learning methods in DTRs focus on the batch (offline) settings, and “development of
5 statistically sound estimation and inference techniques for” the online reinforcement learning (RL) in DTRs “is an
6 important research direction.” We answers this open problem by proposing the first online RL algorithm in DTRs with
7 provably theoretical guarantees. The applications of online RL in health care are motivated by the increasing “use
8 of sophisticated mobile devices” which enables continuous monitoring and intervention on the fly (Chakraborty and
9 Murphy, 2014). For example, a physician could prescribe a set of safe treatments and use online methods to find the
10 optimal combination for a patient with chronic conditions. Broadly speaking, when the combination of the observational
11 data (e.g., $P(v)$) and causal knowledge (causal diagram G) does not ensure the identifiability of the causal effects
12 ($E_\pi[Y]$) in DTRs, the state of art methods would suggest running randomized clinical trials (RCTs) such as SMART
13 (Murphy, 2003, 2005), and solve for the optimal policy using the standard offline methods (e.g., Q-learning). As the
14 reviewer (R2) suggested, running online experiments in the real environment could be dangerous and expensive. Our
15 results are the first adaptive randomization algorithm that could identify the optimal policy in DTRs while achieving
16 near-optimal bounds on the cost of experimentation (regret). For more discussions on adaptive randomization and
17 online RL, see (Gittins, 1979; Rosenberger and Lachin, 2015).

18 **2. High-dimensional State-Action Space (R1, R2)** For experimental studies (e.g., RCTs) in DTRs, issues of sample
19 complexity could arise when the state-action space $|\mathcal{S}||\mathcal{X}|$ is high-dimensional. It was believed in the DTR literature that
20 adaptive randomization procedures (e.g., online RL) could circumvent this issue (Cheung, Chakraborty, and Davidson,
21 2015). Our analysis reveals that this is not the case. In particular, we present first results (Thms. 1-3) analyzing the
22 information complexity of experimental studies in DTRs. We show that a regret bound of $\Omega(\sqrt{|\mathcal{S}||\mathcal{X}|T})$ is inevitable
23 for any randomization procedure, regardless of how sophisticated it might be. This suggests that we should explore
24 other methods to improve the learning convergence. For example, one could exploit the parametric assumptions of the
25 underlying functions (e.g., linear). Our approach takes another route and tries to leverage the abundant, observational
26 data that are available prior to the experimental studies. Specifically, we consider the non-identifiable settings where
27 system dynamics (causal effects $P_{\bar{x}_k}(s_{k+1}|\bar{s}_k)$ and $E_{\bar{x}_k}[Y|\bar{s}_K]$) could not be uniquely determined due to unobserved
28 confounding in canonical DTR models defined in Def. 1 (e.g., Fig. 1(a)). We derive informative bounds about the
29 system dynamics in DTRs from the confounded observational data and incorporate the derived bounds into the online
30 algorithm in an elegant way. We show that this novel approach combining both the online RL and offline bounding
31 could significantly improve the learning performance of online learners for a large family of DTR instances.

32 **3. Identification Conditions (R2)** Our online algorithms (Sec. 2.1) are developed under the conditions of sequential
33 back-door in DTRs (e.g., Fig. 1(b)); while bounding results in Sec. 3.1 are applicable to DTRs with arbitrary unobserved
34 confounding (Fig. 1(a)). Reviewer 2 (R2) seems to be somewhat confused with our identification conditions, and
35 might mistook Fig. 1(a) as the basis of causal assumptions used for online RL methods in Sec. 2.1. As R2 suggested,
36 randomized experiments and online RL are similar in nature (Bareinboim, Forney, and Pearl, 2015). Since each
37 candidate policy π does not take U as input, the sequential backdoor holds in the samples collected by the online RL
38 algorithm in Alg. 1. Causal quantities $P_{\bar{x}_k}(s_{k+1}|\bar{s}_k)$ and $E_{\bar{x}_k}[Y|\bar{s}_K]$ are thus immediately identifiable, and estimable
39 from experimental data. For instance, in Fig. 1(b), S_1, S_2 block all back-door paths from X_1, X_2 to Y . The causal effect
40 $E_{x_1, x_2}[Y|s_1, s_2]$ could be estimated as $E[Y|x_1, x_2, s_1, s_2]$. Given these clarifications, we would like to respectfully
41 ask R2 to re-evaluate our online RL algorithm (Alg. 1) and theoretical regret analysis (Thms. 1-3).

42 **R2:** (1) Regarding the factorization of distribution $P_\pi(\bar{x}_K, \bar{s}_K, y)$ (below Line 113, Page 3), the exogenous U
43 are subsumed in the product of causal quantities $P_{\bar{x}_k}(s_{k+1}|\bar{s}_k)$ and $E_{\bar{x}_k}[Y|\bar{s}_K]$. Specifically, we average U over
44 distribution $P(u)$, and factorize the resultant causal effect $P_{\bar{x}_k}(y, \bar{s}_K)$ following the basic definition of conditional
45 probabilities and exclusion restrictions. (2) Bounds in Thm. 2 is a decreasing function relative to Δ_π and is maximized
46 when Δ_π is the smallest, i.e., π is the second best policy. (3) The consistency axiom suggests that counterfactual
47 probabilities $P(\bar{s}_{k+1_{\bar{x}_k}}|\bar{x}_k) = P(\bar{s}_{k+1}|\bar{x}_k)$. However, Lem. 1 is concerned with gap between causal quantities
48 $P_{\bar{x}_k}(\bar{s}_{k+1}) - P_{\bar{x}_k}(\bar{s}_k)$, which generally do not equate to the bound $P(\bar{s}_{k+1}, \bar{x}_k) - P(\bar{s}_k, \bar{x}_k)$. (4) $\Gamma(s_1)$ is well defined
49 since we define \bar{X}_k as an empty set when $k < 1$. Γ is a function over state-action space for all horizon $k = 1, \dots, K$.
50 (5) In Corol. 1, the upper bound must be no larger than 1 since $E[Y|\bar{s}_K, \bar{x}_K] \leq 1$. (6) Thm. 6 is *stronger* than standard
51 non-identifiability proof since for the constructed DTRs M_1, M_2 , not only their transitional probabilities $P_{\bar{x}_k}(s_{k+1}|\bar{s}_k)$
52 have to be different, but they also need to *match exactly* the lower and upper bounds in Thm. 5. The bounds in Thm. 5
53 are thus tight given the confounded observational data. To the best of our knowledge, we are not aware of any other
54 tightness result regarding DTRs in the literature.

55 **R1, R3:** We really appreciate the reviewers for the helpful suggestions and references. We will incorporate these
56 changes in the camera-ready version of the paper if accepted.