

# Near-Optimal Reinforcement Learning in Dynamic Treatment Regimes

**Junzhe Zhang**

Department of Computer Science  
Columbia University  
junzhez@cs.columbia.edu

**Elias Bareinboim**

Department of Computer Science  
Columbia University  
eb@cs.columbia.edu

## Abstract

A dynamic treatment regime (DTR) consists of a sequence of decision rules, one per stage of intervention, that dictates how to determine the treatment assignment to patients based on evolving treatments and covariates' history. These regimes are particularly effective for managing chronic disorders and is arguably one of the key aspects towards more personalized decision-making. In this paper, we investigate the online reinforcement learning (RL) problem for selecting optimal DTRs provided that observational data is available. We develop the first adaptive algorithm that achieves near-optimal regret in DTRs in online settings, without any access to historical data. We further derive informative bounds on the system dynamics of the underlying DTR from confounded, observational data. Finally, we combine these results and develop a novel RL algorithm that efficiently learns the optimal DTR while leveraging the abundant, yet imperfect confounded observations.

## 1 Introduction

In medical practice, a patient typically has to be treated at multiple stages; the physician repeatedly adapts each treatment, tailored to the patient's time-varying, dynamic state (e.g., level of virus, results of diagnostic tests). Dynamic treatment regimes (DTRs) [20] provide an attractive framework of personalized treatments in longitudinal settings. Operationally, a DTR consists of decision rules that dictate what treatment to provide at each stage, given the patient's evolving conditions and history. These decision rules are alternatively known as adaptive treatment strategies [14, 15, 21, 35, 36] or treatment policies [18, 39, 40]. DTRs offer an effective vehicle for personalized management of chronic conditions, including cancer, diabetes, and mental illnesses [38].

Consider the DTR instance regarding the treatment of alcohol dependence [21, 7], which is graphically represented in Fig. 1a. Based on the condition of alcohol dependant patients ( $S_1$ ), the physician may prescribe a medication or behavioral therapy ( $X_1$ ). Patients are classified as responders or non-responders ( $S_2$ ) based on their level of heavy drinking within the next two months. The physician then must decide whether to continue the initial treatment or switch to an augmented plan combining both medication and behavioral therapy ( $X_2$ ). The unobserved covariate  $U$  summarizes all the unknown factors about the patient. We are interested in the primary outcome  $Y$  that is the percentage of abstinent days over a 12-month period. The treatment policy  $\pi$  in this set-up is a sequence of decision rules  $x_1 \leftarrow \pi_1(s_1), x_2 \leftarrow \pi_2(s_1, s_2, x_1)$  selecting the values of  $X_1, X_2$  based on the history.

Policy learning in a DTR setting is concerned with finding an optimal policy  $\pi$  that maximizes the primary outcome  $Y$ . The main challenge is that since the parameters of the DTR are often unknown, it's not immediate how to directly compute the consequences of executing the policy  $do(\pi)$ , i.e., the expected value  $E_\pi[Y]$ . Most of the current work in the causal inference literature focus on trying to identify this quantity,  $E_\pi[Y]$ , from finite observational data and causal assumptions about the data-generating mechanisms (commonly through causal graphs and potential outcomes). Several criteria

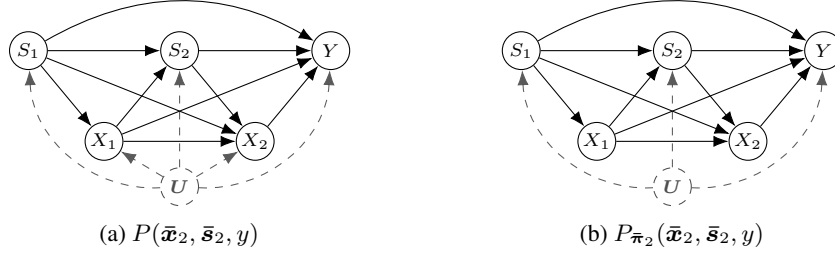


Figure 1: Causal diagrams of (a) a DTR with  $K = 2$  stages of intervention; and (b) a DTR in (a) under sequential interventions  $do(X_1 \sim \pi_1(X_1|S_1), X_2 \sim \pi_2(X_2|S_1, S_2, X_1))$ .

and algorithms have been developed [25, 30, 5]. For instance, a criterion called *sequential backdoor* [26] permits one to determine whether causal effects can be obtained by covariate adjustment. This condition is also referred to as *conditional ignorability* or *unconfoundedness* [29, 20]: there exists no *unobserved confounders* (UCs) that simultaneously affects the treatment at any stage and all the subsequent outcomes given a set of observed covariates. Whenever ignorability holds, a number of efficient estimation procedures exist, including popular methods based on the propensity score [28], inverse probability of treatment weighting [23, 27], and Q-learning [33, 22].

In general, the combination of observational data and causal assumptions does not always lead to point-identification [25, Ch. 3-4]. An alternative is to randomize patients’ treatments at each stage based on the previous decisions and observed outcomes; for instance, one popular strategy is known as the sequential multiple assignment randomized trial (SMART) [21]. By the virtue of randomization, the sequential backdoor condition is entailed. However, in practice, performing a randomized experiment in the actual environment can be extremely costly and undesirable (due to unintended consequences), especially for domains where humans are the main research subjects (e.g., medicine, epidemiology, and psychology). Reinforcement learning (RL) [33] provides a unique opportunity to efficiently learning DTRs due to its nature of balancing exploration and exploitation. A typical RL agent learns by conducting adaptive, sequential experimentation: it repeatedly adjusts the policy that is currently deployed based on the past outcomes. The goal is to learn an optimal policy while minimizing the experimental cost. Efficient RL algorithms have been successfully developed to very general settings such as Markov decision processes (MDPs) [32, 13, 34], where a finite state is statistically sufficient to summarize the treatments and covariates’ history. Variations of this setting include multi-armed bandits [1], partially-observable MDP [11, 3], and factored MDPs [24].

Our focus here is on learning a policy for an unknown DTR while leveraging the observational data. This is a challenging setting for both causal inference and RL. As an example, consider data collected from an unknown behavior policy of the DTR in Fig. 1a (i.e.,  $x_1 \leftarrow f_1(s_1, u)$ ,  $x_2 \leftarrow f_2(s_1, s_2, x_1, u)$ , where both  $U$  and  $\{f_1, f_2\}$  are unobserved), which is materialized in the form of the observational distribution  $P(x_1, x_2, s_1, s_2, y)$  [25, pp. 205]. The existence of the unmeasured confounder  $U$  leads to an immediate violation of the sequential backdoor criterion (e.g., due to the spurious path  $X_1 \leftarrow U \rightarrow Y$ ), which implies that the effect of the policy  $E_{\pi}[Y]$  is not identifiable [25, Ch. 4.4]. On the other hand, existing RL algorithms are not applicable either, which can be seen by noting that DTRs are inherently non-Markovian – in words, the initial treatment  $X_1$  directly affects the outcome  $Y$ . Even though an heuristic approach may be pursued (e.g., Thompson Sampling [37]), and could eventually converge, the same is still not optimal since it’s oblivious to all the observational data.<sup>1</sup> Indeed, it is acknowledged in the literature [8, 9] that the “development of statistically sound estimation and inference techniques” for online RL settings “seem to be another very important research direction”, especially when the increasing use of mobiles devices allows for the possibility of continuous monitoring and just-in-time intervention.

The goal of this paper is to overcome these challenges. We will introduce novel RL strategies capable of optimizing an unknown DTR while efficiently leveraging the imperfect, but large amounts of observational data. In particular, our contributions are as follows: (1) We introduce the first algorithm (UC-DTR (Alg. 1)) that reaches the near-optimal regret bound in the pure DTR setting, without observational data; (2) We derive novel bounds capable of exploiting observational data based on the

<sup>1</sup>Standard off-policy RL methods such as Q-Learning rely on the condition of sequential backdoor, thus not applicable for the confounded observational data. For a more elaborate discussion, see [8, Ch. 3.5]

DTR structure (Thms. 5 and 6), which are provably tight; (3) We develop a novel algorithm (UC<sup>c</sup>-DTR (Alg. 2)) that efficiently incorporates these bounds and accelerates learning in the online setting. Our results are validated on randomly generated DTRs and multi-stage clinical trials on cancer treatment.

## 1.1 Preliminaries

In this section, we introduce the basic notation and definitions used throughout the paper. We use capital letters to denote variables ( $X$ ) and small letters for their values ( $x$ ). Let  $\mathcal{X}$  represent the domain of  $X$  and  $|\mathcal{X}|$  its dimension. We consistently use the abbreviation  $P(x)$  to represent the probabilities  $P(X = x)$ .  $\bar{X}_k$  stands for a sequence  $\{X_1, \dots, X_k\}$  ( $\emptyset$  if  $k < 1$ ), and  $\bar{\mathcal{X}}_k$  represents its domain, i.e.,  $\mathcal{X}_1 \times \dots \times \mathcal{X}_k$ . Further, we denote by  $I_{\{\cdot\}}$  the indicator function.

The basic semantical framework of our analysis rests on *structural causal models* (SCM) [25, Ch. 7]. A SCM  $M$  is a tuple  $\langle U, V, F, P(u) \rangle$  where  $U$  is a set of exogenous (unobserved) variables and  $V$  is a set of endogenous (observed) variables.  $F$  is a set of structural functions where  $f_i \in F$  decides the values of  $V_i \in V$  taking as argument a combination of other endogenous and exogenous variables (i.e.,  $V_i \leftarrow f_i(PA_i, U_i)$ ,  $PA_i \subseteq V, U_i \subseteq U$ ). The values of  $U$  are drawn from the distribution  $P(u)$ , and induce an observational distribution  $P(v)$  [25, pp. 205]. Each SCM is associated with a causal diagram in the form of a directed acyclic graph  $G$ , where nodes represent endogenous variables, dashed nodes exogenous variables, and arrows stand for functional relations (e.g., see Fig. 1).

An intervention on a set of endogenous variables  $X$ , denoted by  $do(x)$ , is an operation where values of  $X$  are set to constants  $x$ , regardless of how they were ordinarily determined (through the functions  $\{f_X : \forall X \in \mathbf{X}\}$ ). For a SCM  $M$ , let  $M_x$  be a sub-model of  $M$  induced by intervention  $do(x)$ . The interventional distribution  $P_x(y)$  induced by  $do(x)$  is the distribution over variables  $Y$  in the sub-model  $M_x$ . For a more detailed discussion of SCMs, we refer readers to [25, Ch. 7].

## 2 Optimizing Dynamic Treatment Regimes

In this section, we will formalize the problem of online optimization in DTRs with confounded observations and provide an efficient solution. We start by defining DTRs in the structural semantics.

**Definition 1** (Dynamic Treatment Regime [20]). A dynamic treatment regime (DTR) is a SCM  $\langle U, V, F, P(u) \rangle$  where the endogenous variables  $V = \{\bar{X}_K, \bar{S}_K, Y\}$ ;  $K \in \mathbb{N}^+$  is the total stages of interventions. For stage  $k = 1, \dots, K$ : (1)  $X_k$  is a finite decision decided by a behavior policy  $x_k \leftarrow f_k(\bar{s}_k, \bar{x}_{k-1}, u)$ ; (2)  $S_k$  is a finite state decided by a transition function  $s_k \leftarrow \tau_k(\bar{x}_{k-1}, \bar{s}_{k-1}, u)$ .  $Y$  is the primary outcome at the final state  $K$ , decided by a reward function  $y \leftarrow r(\bar{x}_K, \bar{s}_K, u)$  bounded in  $[0, 1]$ . Values of exogenous variables  $U$  are drawn from the distribution  $P(u)$ .

A DTR  $M^*$  induces an observational distribution  $P(\bar{x}_K, \bar{s}_K, y)$ . Fig. 1a shows the causal diagram of a DTR with  $K = 2$  stages of interventions. A policy  $\pi$  for a DTR is a sequence of decision rules  $\bar{\pi}_K$ , where each  $\pi_k(x_k | \bar{s}_k, \bar{x}_{k-1})$  is a function mapping from the domain of histories  $\bar{S}_k, \bar{X}_{k-1}$  up to stage  $k$  to a distribution over decision  $X_k$ . A policy is called *deterministic* if the above mappings are from histories  $\bar{S}_k, \bar{X}_{k-1}$  to the domain of decision  $X_k$ , i.e.,  $x_k \leftarrow \pi_k(\bar{s}_k, \bar{x}_{k-1})$ . The collection of possible policies, depending on the domains of the history and decision, define a policy space  $\Pi$ .

A policy  $\pi$  defines a sequence of stochastic interventions  $do(X_1 \sim \pi_1(X_1 | \bar{S}_1), \dots, X_K \sim \pi_K(X_K | \bar{S}_K, \bar{X}_{K-1}))$ , which induce an interventional distribution over variables  $\bar{X}_K, \bar{S}_K, Y$ , i.e.:

$$P_\pi(\bar{x}_K, \bar{s}_K, y) = P_{\bar{x}_K}(y | \bar{s}_K) \prod_{k=0}^{K-1} P_{\bar{x}_k}(s_{k+1} | \bar{s}_k) \pi_{k+1}(x_{k+1} | \bar{s}_{k+1}, \bar{x}_k), \quad (1)$$

where  $P_{\bar{x}_k}(s_{k+1} | \bar{s}_k)$  is the transition distribution at stage  $k$  and  $P_{\bar{x}_K}(y | \bar{s}_K)$  is the reward distribution over the primary outcome. Fig. 1b describes a DTR under  $K = 2$  stages of interventions  $do(X_2 \sim \pi_1(X_1 | S_1), X_2 \sim \pi_2(X_2 | S_1, S_2, X_1))$ . The expected cumulative reward of a policy  $\pi$  in a DTR  $M^*$  is given by  $V_\pi(M^*) = E_\pi[Y]$ . We are searching for an optimal policy  $\pi^*$  that maximizes the cumulative reward, i.e.,  $\pi^* = \arg \max_{\pi \in \Pi} V_\pi(M^*)$ . It is a well-known fact in decision theory that no stochastic policy can improve on the utility of the best deterministic policy (see, e.g., [17, Lem. 2.1]). Thus, in what follows, we will usually consider the policy space  $\Pi$  to be deterministic.

Our goal is to optimize an unknown DTR  $M^*$  based solely on the domains  $\mathcal{S} = \bar{\mathcal{S}}_K, \mathcal{X} = \bar{\mathcal{X}}_K$  and the observational distribution  $P(\bar{x}_K, \bar{s}_K, y)$  (i.e., both  $F, P(u)$  are unknown). The agent (e.g., a

---

**Algorithm 1:** UC-DTR

---

**Input:** failure tolerance  $\delta \in (0, 1)$ .

- 1: **for all** episodes  $t = 1, 2, \dots$  **do**
- 2: Define event counts  $N^t(\bar{s}_k, \bar{x}_k)$  and  $N^t(\bar{s}_k, \bar{x}_{k-1})$  for horizon  $k = 1, \dots, K$  prior to episode  $t$  as, respectively,  $\sum_{i=1}^{t-1} I_{\bar{s}_k^i = \bar{s}_k, \bar{x}_k^i = \bar{x}_k}$  and  $\sum_{i=1}^{t-1} I_{\bar{s}_k^i = \bar{s}_k, \bar{x}_{k-1}^i = \bar{x}_{k-1}}$ . Further, define reward counts  $R^t(\bar{s}_K, \bar{x}_K)$  prior to episode  $t$  as  $\sum_{i=1}^{t-1} Y^i I_{\bar{s}_K^i = \bar{s}_K, \bar{x}_K^i = \bar{x}_K}$ .
- 3: Compute estimates  $\hat{P}_{\bar{x}_k}^t(s_{k+1}|\bar{s}_k)$  and  $\hat{E}_{\bar{x}_K}^t[Y|\bar{s}_K]$  as

$$\hat{P}_{\bar{x}_k}^t(s_{k+1}|\bar{s}_k) = \frac{N^t(\bar{s}_{k+1}, \bar{x}_k)}{\max\{1, N^t(\bar{s}_k, \bar{x}_k)\}}, \quad \hat{E}_{\bar{x}_K}^t[Y|\bar{s}_K] = \frac{R^t(\bar{s}_K, \bar{x}_K)}{\max\{1, N^t(\bar{s}_k, \bar{x}_k)\}}.$$

- 4: Let  $\mathcal{M}_t$  denote a set of DTRs such that for any  $M \in \mathcal{M}_t$ , its transition probabilities  $P_{\bar{x}_k}(s_{k+1}|\bar{s}_k)$  and reward  $E_{\bar{x}_K}[Y|\bar{s}_K]$  are close to estimates  $\hat{P}_{\bar{x}_k}^t(s_{k+1}|\bar{s}_k)$ ,  $\hat{E}_{\bar{x}_K}^t[Y|\bar{s}_K]$ , i.e.,

$$\left\| P_{\bar{x}_k}(\cdot|\bar{s}_k) - \hat{P}_{\bar{x}_k}^t(\cdot|\bar{s}_k) \right\|_1 \leq \sqrt{\frac{6|\mathcal{S}_{k+1}| \log(2K|\mathcal{S}_k||\mathcal{X}_k|t/\delta)}{\max\{1, N^t(\bar{s}_k, \bar{x}_k)\}}}, \quad (2)$$

$$\left| E_{\bar{x}_K}[Y|\bar{s}_K] - \hat{E}_{\bar{x}_K}^t[Y|\bar{s}_K] \right| \leq \sqrt{\frac{2 \log(2K|\mathcal{S}||\mathcal{X}|t/\delta)}{\max\{1, N^t(\bar{s}_K, \bar{x}_K)\}}}. \quad (3)$$

- 5: Find the optimal policy  $\pi_t$  of an optimistic DTR  $M_t \in \mathcal{M}_t$  such that

$$V_{\pi_t}(M_t) = \max_{\pi \in \Pi, M \in \mathcal{M}_t} V_{\pi}(M) \quad (4)$$

- 6: Execute policy  $\pi_t$  for episode  $t$  and observe the samples  $\bar{S}_K^t, \bar{X}_K^t, Y^t$ .
  - 7: **end for**
- 

physician) learns through repeated experiments of episodes  $t = 1, \dots, T$ . Each episode  $t$  contains a complete DTR process: at stage  $k$ , the agent observes the state  $S_k^t$ , performs an intervention  $do(X_k^t)$  and moves to the state  $S_{k+1}^t$ ; the primary outcome  $Y^t$  is received at the final stage  $K$ . The cumulative regret up to episode  $T$  is defined as  $R(T) = \sum_{t=1}^T (V_{\pi^*}(M^*) - Y^t)$ , i.e., the loss due to the fact that the agent does not always pick the optimal policy  $\pi^*$ . We will assess and compare algorithms in terms of their regret  $R(T)$ . A desirable asymptotic property is to have  $\lim_{T \rightarrow \infty} E[R(T)]/T = 0$ , meaning that the agent eventually converges and finds the optimal policy  $\pi^*$ .

## 2.1 The UC-DTR Algorithm

We now introduce a new RL algorithm for optimizing an unknown DTR, which we call UC-DTR. We will later prove that UC-DTR achieves near-optimal bound on the total regret given only the knowledge of the domains  $\mathcal{S}$  and  $\mathcal{X}$ . Like many other online RL algorithms [1, 13, 24], UC-DTR follows the principle of *optimism under uncertainty* to balance exploration and exploitation. The algorithm generally works in phases of model learning, optimistic planning, and strategy execution.

The details of UC-DTR procedure can be found in Alg. 1. The algorithm proceeds in episodes and computes a new strategy  $\pi_t$  from samples  $\{\bar{S}_K^i, \bar{X}_K^i, Y^i\}_{i=1}^{t-1}$  collected so far at the beginning of each episode  $t$ . Specifically, UC-DTR computes in Steps 1-3, the empirical estimates  $\hat{E}_{\bar{x}_K}^t[Y|\bar{s}_K]$  of the expected reward  $E_{\bar{x}_K}[Y|\bar{s}_K]$ , and  $\hat{P}_{\bar{x}_k}^t(s_{k+1}|\bar{s}_k)$  of the transitional probabilities  $P_{\bar{x}_k}(s_{k+1}|\bar{s}_k)$  from experimental samples collected prior to episode  $t$ . In Step 4, a set  $\mathcal{M}_t$  of plausible DTRs is defined in terms of confidence region around the the empirical estimates  $\hat{E}_{\bar{x}_K}^t[Y|\bar{s}_K]$  and  $\hat{P}_{\bar{x}_k}^t(s_{k+1}|\bar{s}_k)$ . This guarantees that the true DTR  $M^*$  is in the set  $\mathcal{M}_t$  with high probability. In Step 5, UC-DTR computes the optimal policy  $\pi_t$  of the most optimistic instance  $M_t$  in the family of DTRs  $\mathcal{M}_t$  that induces the maximal optimal expected reward. This policy  $\pi_t$  is executed throughout episode  $t$  and new samples  $\bar{S}_K^t, \bar{X}_K^t, Y^t$  are collected (Step 6).

**Finding Optimistic DTRs** The Step 5 of UC-DTR tries to find an optimal policy  $\pi_t$  for an optimistic DTR  $M_t$ . While the Bellman equation [6] allows one to optimize a fixed DTR, we need to find a DTR  $M_t$  that gives the maximal optimal reward among all plausible DTRs in  $\mathcal{M}_t$  given by Eq. (3).

We now introduce a method that extends standard dynamic programming planners [6] to solve this problem. We first combine all DTRs in  $\mathcal{M}_t$  to get an extended DTR  $M_+$ . That is, we consider a DTR  $M_+$  with continuous decision space  $\bar{\mathcal{X}}^+ = \bar{\mathcal{X}}_K^+$ , where for each horizon  $k$ , each action  $\bar{x}_k \in \bar{\mathcal{X}}_k^+$ , each admissible transition probabilities  $P_{\bar{x}_k}(s_{k+1}|\bar{s}_k)$  according to Eq. (2), there is an action in  $\bar{\mathcal{X}}_k^+$  inducing the same probabilities  $P_{\bar{x}_k}(s_{k+1}|\bar{s}_k)$ . Similar arguments also apply to the expected reward  $E_{\bar{x}_K}[Y|\bar{s}_K]$ . Then, for each policy  $\pi_+$  on  $M_+$ , there is an DTR  $M_t \in \mathcal{M}_t$  and a policy  $\pi_t \in \Pi$  such that policies  $\pi_+$  and  $\pi_t$  induces the same transition probabilities on the respective DTR, and vice versa. Thus, solving the optimization problem in Eq. (4) is equivalent to finding an optimal policy  $\pi_+^*$  on the extended DTR  $M_+$ . Let  $V^*(\bar{s}_k, \bar{x}_{k-1})$  denote the optimal value  $E_{\pi_+^*}[Y|\bar{s}_k, \bar{x}_{k-1}]$  in  $M_+$ . The Bellman equation on  $M_+$  for  $k = 1, \dots, K-1$  is defined as follows:

$$V^*(\bar{s}_k, \bar{x}_{k-1}) = \max_{x_k} \left\{ \max_{P_{\bar{x}_k}(\cdot|\bar{s}_k) \in \mathcal{P}_k} \left\{ \sum_{s_{k+1}} V^*(\bar{s}_{k+1}, \bar{x}_k) P_{\bar{x}_k}(s_{k+1}|\bar{s}_k) \right\} \right\}, \quad (5)$$

and  $V^*(\bar{s}_K, \bar{x}_{K-1}) = \max_{x_K} \max_{E_{\bar{x}_K}[Y|\bar{s}_K] \in \mathcal{R}} E_{\bar{x}_K}[Y|\bar{s}_K],$

where  $\mathcal{R}$  and  $\mathcal{P}_k$  are the convex polytope of parameters  $E_{\bar{x}_K}[Y|\bar{s}_K]$  and  $P_{\bar{x}_k}(s_{k+1}|\bar{s}_k)$  defined in Eqs. (2) and (3), respectively. The inner maximum in Eq. (5) is a linear program (LP) over the convex polytope  $\mathcal{P}_k$  (or  $\mathcal{R}$ ), which is solvable using standard LP algorithms.

## 2.2 Theoretical Analysis

We now analyze the asymptotic behavior of UC-DTR that will lead to a better understanding of its theoretical guarantees. Given space constraints, all proofs are provided in Appendix I. The following proposition shows that the cumulative regret of UC-DTR after  $T$  steps is at most  $\tilde{O}(K\sqrt{|\mathcal{S}||\mathcal{X}|T})^2$ .

**Theorem 1.** Fix a  $\delta \in (0, 1)$ . With probability (w.p.) of at least  $1 - \delta$ , it holds for any  $T > 1$ , the regret of UC-DTR with parameter  $\delta$  is bounded by

$$R(T) \leq 12K\sqrt{|\mathcal{S}||\mathcal{X}|T \log(2K|\mathcal{S}||\mathcal{X}|T/\delta)} + 4K\sqrt{T \log(2T/\delta)}. \quad (6)$$

It is also possible to obtain the instance-dependent bound on the expected regret. Let  $\Pi^-$  denote a set of sub-optimal policies  $\{\pi \in \Pi : V_\pi(M^*) < V_{\pi^*}(M^*)\}$ . For any  $\pi \in \Pi^-$ , let its gap in expected reward between the optimal policy  $\pi^*$  be  $\Delta_\pi = V_{\pi^*}(M^*) - V_\pi(M^*)$ . We next derive the gap-dependent logarithmic bound on the expected regret of UC-DTR after  $T$  steps.

**Theorem 2.** For any  $T \geq 1$ , with parameter  $\delta = \frac{1}{T}$ , the expected regret of UC-DTR is bounded by

$$E[R(T)] \leq \max_{\pi \in \Pi^-} \left\{ \frac{33^2 K^2 |\mathcal{S}||\mathcal{X}| \log(T)}{\Delta_\pi} + \frac{32}{\Delta_\pi^3} + \frac{4}{\Delta_\pi} \right\} + 1. \quad (7)$$

Since Eq. (7) is a decreasing function relative to the gap  $\Delta_\pi$ , the maximum of the regret in Thm. 2 is achieved with the second best policy  $\pi^- = \arg \min_{\pi \in \Pi^-} \Delta_\pi$ . We also provide a corresponding lower bound on the expected regret of any experimental algorithm.

**Theorem 3.** For any algorithm  $\mathcal{A}$ , any natural numbers  $K \geq 1$ , and  $|\mathcal{S}^k| \geq 2, |\mathcal{X}^k| \geq 2$  for any  $k \in \{1, \dots, K\}$ , there is a DTR  $M$  with horizon  $K$ , state domains  $\mathcal{S}$  and action domains  $\mathcal{X}$ , such that the expected regret of  $\mathcal{A}$  after  $T \geq |\mathcal{S}||\mathcal{X}|$  episodes is at least

$$E[R(T)] \geq 0.05\sqrt{|\mathcal{S}||\mathcal{X}|T} \quad (8)$$

Thm. 3 implies that for any DTR instance, the cumulative regret of  $\Omega(\sqrt{|\mathcal{S}||\mathcal{X}|T})$  is inevitable. The regret upper bound  $\tilde{O}(K\sqrt{|\mathcal{S}||\mathcal{X}|T})$  in Thm. 1 is close to the lower bound  $\Omega(\sqrt{|\mathcal{S}||\mathcal{X}|T})$  in Thm. 3, which means that UC-DTR is near-optimal provided with only the domains of state  $\mathcal{S}$  and actions  $\mathcal{X}$ .

<sup>2</sup> $\tilde{O}(\cdot)$  is similar to  $\mathcal{O}(\cdot)$  but ignores log-terms, i.e.,  $f = \tilde{O}(g)$  if and only if  $\exists k, f = \mathcal{O}(g \log^k(g))$ .

### 3 Learning from Confounded Observations

The results presented so far (Thms. 1 to 3) establish the dimension of the state-action domain  $|\mathcal{S}||\mathcal{X}|$  as the an important parameter for the information complexity of online learning in DTRs. When domains  $\mathcal{S} \times \mathcal{X}$  are high-dimensional, the cumulative regret will be significant for any online algorithm, no matter how sophisticated it might be. This observation suggests that we should explore other reasonable assumptions to address the issues of high-dimensional domains.

A natural approach is to utilize the abundant observational data, which could be obtained by passively observing other agents behaving in the environment. Despite all its power, the UC-DTR algorithm does not make use of any knowledge in the the observational distribution  $P(\bar{s}_K, \bar{x}_K, y)$ . For the remainder of this paper, we will present and study an efficient procedure to incorporate observational samples of  $P(\bar{s}_K, \bar{x}_K, y)$ , so that the performance of online learners could be improved.

When states  $\bar{S}_K$  satisfy the *sequential backdoor* criterion [26] with respect to treatments  $\bar{X}_K$  and the primary outcome  $Y$ , one could identify the transition probabilities  $P_{\bar{x}_k}(s_{k+1}|\bar{s}_k)$  and expected reward  $E_{\bar{x}_K}[Y|\bar{s}_k]$  from  $P(\bar{s}_K, \bar{x}_K, y)$ . The optimal policy is thus solvable using the standard off-policy learning methods such as Q-learning [33, 22]. However, issues of non-identifiability arise in the general settings where the sequential backdoor does not hold (e.g., see Fig. 1a).

**Theorem 4.** *Given  $P(\bar{s}_K, \bar{x}_K, y) > 0$ , there exists DTRs  $M_1, M_2$  such that  $P^{M_1}(\bar{s}_K, \bar{x}_K, y) = P^{M_2}(\bar{s}_K, \bar{x}_K, y) = P(\bar{s}_K, \bar{x}_K, y)$  while  $P_{\bar{x}_K}^{M_1}(\bar{s}_K, y) \neq P_{\bar{x}_K}^{M_2}(\bar{s}_K, y)$ .*

Thm. 4 is stronger than the standard non-identifiability results (e.g., [16, Thm. 1]). It shows that given *any* observational distribution  $P(\bar{s}_K, \bar{x}_K, y)$ , one to construct two DTRs both compatible with  $P(\bar{s}_K, \bar{x}_K, y)$ , but disagrees in the interventional probabilities  $P_{\bar{x}_K}(\bar{s}_K, y)$ .

#### 3.1 Bounds and Partial Identification in DTRs

In this section, we consider a partial identification task in DTRs which bounds parameters of  $P_{\bar{x}_k}(s_{k+1}|\bar{s}_k)$  and  $E_{\bar{x}_K}[Y|\bar{s}_k]$  from the observational distribution  $P(\bar{s}_K, \bar{x}_K, y)$ . Our first result shows that the gap between causal quantities  $P_{\bar{x}_k}(s_{k+1})$  and  $P_{\bar{x}_k}(\bar{s}_k)$  in a DTR is bounded by the gap between the corresponding observational distributions  $P(\bar{s}_{k+1}, \bar{x}_k)$  and  $P(\bar{s}_k, \bar{x}_k)$ .

**Lemma 1.** *For a DTR, given  $P(\bar{s}_K, \bar{x}_K, y)$ , for any  $k = 1, \dots, K-1$ ,*

$$P_{\bar{x}_k}(\bar{s}_{k+1}) - P_{\bar{x}_k}(\bar{s}_k) \leq P(\bar{s}_{k+1}, \bar{x}_k) - P(\bar{s}_k, \bar{x}_k). \quad (9)$$

Lem. 1 allows one to derive informative bounds of transition probabilities  $P_{\bar{x}_k}(s_{k+1}|\bar{s}_k)$  in a DTR, which are consistently estimable from the observational data  $P(\bar{s}_K, \bar{x}_K)$ .

**Theorem 5.** *For a DTR, given  $P(\bar{s}_K, \bar{x}_K, y) > 0$ , for any  $k = 1, \dots, K-1$ ,*

$$\frac{P(\bar{s}_{k+1}, \bar{x}_k)}{\Gamma(\bar{s}_k, \bar{x}_{k-1})} \leq P_{\bar{x}_k}(s_{k+1}|\bar{s}_k) \leq \frac{\Gamma(\bar{s}_{k+1}, \bar{x}_k)}{\Gamma(\bar{s}_k, \bar{x}_{k-1})}, \quad (10)$$

where  $\Gamma(\bar{s}_{k+1}, \bar{x}_k) = P(\bar{s}_{k+1}, \bar{x}_k) - P(\bar{s}_k, \bar{x}_k) + \Gamma(\bar{s}_k, \bar{x}_{k-1})$  and  $\Gamma(s_1) = P(s_1)$ .

Bounds in Thm. 5 exploit the sequential functional relationships among states and treatments in the underlying DTR, which improve over the best-known bounds reported in [19, 4, 41]. Let  $[a_{\bar{x}_k, \bar{s}_k}(s_{k+1}), b_{\bar{x}_k, \bar{s}_k}(s_{k+1})]$  denote the bound over  $P_{\bar{x}_k}(s_{k+1}|\bar{s}_k)$  given by Eq. (10). We next show that  $P_{\bar{x}_k}(s_{k+1}|\bar{s}_k) \in [a_{\bar{x}_k, \bar{s}_k}(s_{k+1}), b_{\bar{x}_k, \bar{s}_k}(s_{k+1})]$  is indeed optimal without additional assumption.

**Theorem 6.** *Given  $P(\bar{s}_K, \bar{x}_K, y) > 0$ , for any  $k \in \{1, \dots, K-1\}$ , there exists DTRs  $M_1, M_2$  such that  $P^{M_1}(\bar{s}_K, \bar{x}_K, y) = P^{M_2}(\bar{s}_K, \bar{x}_K, y) = P(\bar{s}_K, \bar{x}_K, y)$  while  $P_{\bar{x}_k}^{M_1}(s_{k+1}|\bar{s}_k) = a_{\bar{x}_k, \bar{s}_k}(s_{k+1})$ ,  $P_{\bar{x}_k}^{M_2}(s_{k+1}|\bar{s}_k) = b_{\bar{x}_k, \bar{s}_k}(s_{k+1})$ .*

Thm. 6 ensures the optimality of Thm. 5. Suppose there exists a bound  $[a'_{\bar{x}_k, \bar{s}_k}(s_{k+1}), b'_{\bar{x}_k, \bar{s}_k}(s_{k+1})]$  strictly contained in  $[a_{\bar{x}_k, \bar{s}_k}(s_{k+1}), b_{\bar{x}_k, \bar{s}_k}(s_{k+1})]$ . By Thm. 6, we could always find DTRs  $M_1, M_2$  that are compatible with the observational data  $P(\bar{s}_K, \bar{x}_K, y)$  while their transition probabilities  $P_{\bar{x}_k}(s_{k+1}|\bar{s}_k)$  lie outside of the bound  $[a'_{\bar{x}_k, \bar{s}_k}(s_{k+1}), b'_{\bar{x}_k, \bar{s}_k}(s_{k+1})]$ , which is a contradiction.

As a corollary, one could apply methods of Lem. 1 and Thm. 5 to bound expected rewards  $E_{\bar{x}_K}[Y|\bar{s}_k]$  from  $P(\bar{s}_K, \bar{x}_K, y)$ . The optimality of the derived bounds follows immediately after Thm. 6.

---

**Algorithm 2:** Causal UC-DTR (UC<sup>c</sup>-DTR)

---

**Input:** failure tolerance  $\delta \in (0, 1)$ , causal bounds  $\mathcal{C}$ .

- 1: Let  $\mathcal{M}^c$  denote a set of DTRs compatible with causal bounds  $\mathcal{C}$ , i.e., for any  $M \in \mathcal{M}^c$ , its causal quantities  $P_{\bar{x}_k}(s_{k+1}|\bar{s}_k)$  and  $E_{\bar{x}_K}[Y|\bar{s}_K]$  satisfy Eq. (13) and Eq. (14) respectively.
- 2: **for all** episodes  $t = 1, 2, \dots$  **do**
- 3:   Execute Steps 2-4 of UC-DTR (Alg. 1).
- 4:   Find the optimal policy  $\pi_t$  of an optimistic DTR  $M_t$  in  $\mathcal{M}_t^c = \mathcal{M}_t \cap \mathcal{M}^c$  such that

$$V_{\pi_t}(M_t) = \max_{\pi \in \Pi, M \in \mathcal{M}_t^c} V_{\pi}(M) \quad (12)$$

- 5:   Execute policy  $\pi_t$  for episode  $t$  and observe the samples  $\bar{S}_K^t, \bar{X}_K^t, Y^t$ .
  - 6: **end for**
- 

**Corollary 1.** For a DTR, given  $P(\bar{s}_K, \bar{x}_K, y) > 0$ ,

$$\frac{E[Y|\bar{s}_K, \bar{x}_K]P(\bar{s}_K, \bar{x}_K)}{\Gamma(\bar{s}_K, \bar{x}_{K-1})} \leq E_{\bar{x}_K}[Y|\bar{s}_K] \leq 1 - \frac{(1 - E[Y|\bar{s}_K, \bar{x}_K])P(\bar{s}_K, \bar{x}_K)}{\Gamma(\bar{s}_K, \bar{x}_{K-1})}. \quad (11)$$

Since  $E[Y|\bar{s}_K, \bar{x}_K] \in [0, 1]$ , the bounds in Eq. (11) are contained in  $[0, 1]$  and are thus informative. The bounds developed so far are functions of the observational distribution  $P(\bar{s}_K, \bar{x}_K, y)$  which is identifiable by the sampling process, and so generally can be estimated consistently. Specifically, we estimate the bounds in Thm. 5 and Corol. 1 by the corresponding sample mean estimates. Standard results of large-deviation theory are thus applicable to control the uncertainties due to finite samples.

### 3.2 The Causal UC-DTR Algorithm

Our goal in this section is to introduce a simple, yet principled approach for leveraging the new-found bounds defined in Thm. 5 and Corol. 1, hopefully improving the performance of UC-DTR procedure.

For  $k = 1, \dots, K-1$ , let  $\mathcal{C}_k$  denote a set of bounds over transition probabilities  $P_{\bar{x}_k}(s_{k+1}|\bar{s}_k)$ , i.e.,

$$\mathcal{C}_k = \left\{ \forall \bar{s}_{k+1}, \bar{x}_k : P_{\bar{x}_k}(s_{k+1}|\bar{s}_k) \in [a_{\bar{x}_k, \bar{s}_k}(s_{k+1}), b_{\bar{x}_k, \bar{s}_k}(s_{k+1})] \right\}. \quad (13)$$

Similarly, let  $\mathcal{C}_K$  denote a set of bounds over the conditional expected reward  $E_{\bar{x}_K}[Y|\bar{s}_K]$ , i.e.,

$$\mathcal{C}_K = \left\{ \forall \bar{s}_K, \bar{x}_K : E_{\bar{x}_K}[Y|\bar{s}_K] \in [a_{\bar{x}_K, \bar{s}_K}, b_{\bar{x}_K, \bar{s}_K}] \right\}. \quad (14)$$

We denote by  $\mathcal{C}$  a set of bounds  $\{\mathcal{C}_1, \dots, \mathcal{C}_K\}$  on the system dynamics of the DTR, called *causal bounds*. Our procedure Causal UC-DTR (for short, UC<sup>c</sup>-DTR) is summarized in Alg. 2. UC<sup>c</sup>-DTR is similar to the original UC-DTR but exploits causal bounds  $\mathcal{C}$ . It maintains a set of possible DTRs  $\mathcal{M}^c$  compatible with the causal bounds  $\mathcal{C}$  (Step 1). Before each episode  $t$ , it computes the optimal policy  $\pi_t$  of an optimistic DTRs  $M_t$  in set  $\mathcal{M}_t^c = \mathcal{M}_t \cap \mathcal{M}^c$  (Step 3). Similar to UC-DTR,  $\pi_t$  could be obtained by solving LPs defined in Eq. (5) subject to additional causal constraints Eqs. (13) and (14).

We next analyze asymptotic properties of UC<sup>c</sup>-DTR, showing that it consistently outperforms UC-DTR. Let  $\|\mathcal{C}_k\|_1$  denote the maximal L1 norm of any parameter in  $\mathcal{C}_k$ , i.e., for any  $k = 1, \dots, K-1$ ,

$$\|\mathcal{C}_k\|_1 = \max_{\bar{x}_k, \bar{s}_k} \sum_{s_{k+1}} |a_{\bar{x}_k, \bar{s}_k}(s_{k+1}) - b_{\bar{x}_k, \bar{s}_k}(s_{k+1})|, \quad \text{and} \quad \|\mathcal{C}_K\|_1 = \max_{\bar{x}_K, \bar{s}_K} |a_{\bar{x}_K, \bar{s}_K} - b_{\bar{x}_K, \bar{s}_K}|.$$

Further, let  $\|\mathcal{C}\|_1 = \sum_{k=1}^K \|\mathcal{C}_k\|_1$ . The total regret of UC<sup>c</sup>-DTR after  $T$  steps is bounded as follows.

**Theorem 7.** Fix a  $\delta \in (0, 1)$ . With probability of at least  $1 - \delta$ , it holds for any  $T > 1$ , the regret of UC<sup>c</sup>-DTR with parameter  $\delta$  and causal bounds  $\mathcal{C}$  is bounded by

$$R(T) \leq \min \left\{ 12K \sqrt{|\mathcal{S}||\mathcal{X}|T \log(2K|\mathcal{S}||\mathcal{X}|T/\delta)}, \|\mathcal{C}\|_1 T \right\} + 4K \sqrt{T \log(2T/\delta)}. \quad (15)$$

It is immediate from Thm. 7 that the regret bound in Eq. (15) is smaller than the bound given by Eq. (6) if  $T < 12^2 |\mathcal{S}||\mathcal{X}| \log(2K|\mathcal{S}||\mathcal{X}|T/\delta) / \|\mathcal{C}\|_1^2$ . This means that UC<sup>c</sup>-DTR has a head start over UC-DTR when the causal bounds  $\mathcal{C}$  are informative, i.e., the dimension  $\|\mathcal{C}\|_1$  is small.

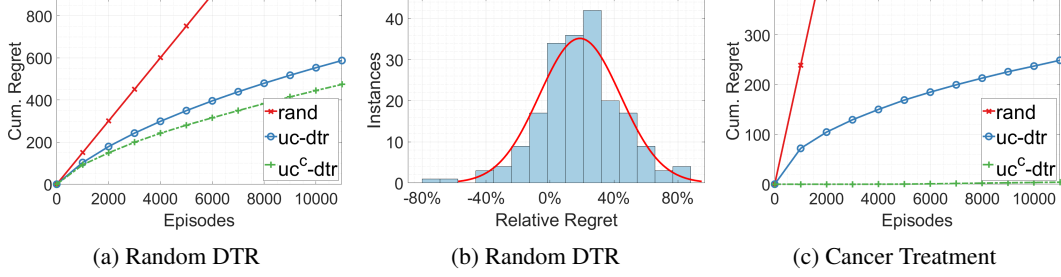


Figure 2: Simulations comparing online learners that are randomized (*rand*), adaptive (*uc-dtr*) and causally enhanced (*uc<sup>c</sup>-dtr*). Graphs are rendered in high resolution and can be zoomed in.

We could also witness the improvements of causal bounds on the total expected regret. Let  $\Pi_{\mathcal{C}}^-$  be the set of sub-optimal policies that their maximal expected rewards over instances in  $\mathcal{M}^c$  are no less than the true optimal value  $V_{\pi^*}(M^*)$ , i.e.,  $\Pi_{\mathcal{C}}^- = \{\pi \in \Pi^- : \max_{M \in \mathcal{M}^c} V_{\pi}(M) \geq V_{\pi^*}(M^*)\}$ . The following is the instance-dependent bound on the total regret of UC<sup>c</sup>-DTR after  $T$  steps.

**Theorem 8.** For any  $T \geq 1$ , with parameter  $\delta = \frac{1}{T}$  and causal bounds  $\mathcal{C}$ , the expected regret of UC<sup>c</sup>-DTR is bounded by

$$E[R(T)] \leq \max_{\pi \in \Pi_{\mathcal{C}}^-} \left\{ \frac{33^2 K^2 |\mathcal{S}| |\mathcal{X}| \log(T)}{\Delta_{\pi}} + \frac{32}{\Delta_{\pi}^3} + \frac{4}{\Delta_{\pi}} \right\} + 1. \quad (16)$$

Since  $\Pi_{\mathcal{C}}^- \subseteq \Pi^-$ , it follows that the regret bound in Thm. 8 is small than or equal to Eq. (7), i.e., UC<sup>c</sup>-DTR consistently dominates UC-DTR in terms of the performance. For instance, in a multi-armed bandit model (i.e., 1-stage DTR with  $S_1 = \emptyset$ ) with optimal reward  $\mu^*$ , the regret of UC<sup>c</sup>-DTR is  $\mathcal{O}(|\mathcal{X}| \log(T)/\Delta_x)$  where  $\Delta_x$  is the smallest gap among sub-optimal arms  $x$  satisfying  $b_x \geq \mu^*$ .

## 4 Experiments

We demonstrate our algorithms on several dynamic treatment regimes, including randomly generated DTRs, and the survival model in the context of multi-stage cancer treatment. We found that our algorithms could efficiently found the optimal policy; the observational data typically improve the convergence rate of online RL learners despite the confounding bias.

In all experiments, we test sequentially randomized trials (*rand*), UC-DTR algorithm (*uc-dtr*) and the causal UC-DTR (*uc<sup>c</sup>-dtr*) with causal bounds derived from  $1 \times 10^5$  confounded observational samples. Each experiment lasts for  $T = 1.1 \times 10^4$  episodes. The parameter  $\delta = \frac{1}{KT}$  for *uc-dtr* and *uc<sup>c</sup>-dtr* where  $K$  is the total stages of interventions. For all algorithms, we measure their cumulative regret over 200 repetitions. We refer readers to Appendix II for the more details on the experimental set-up.

**Random DTRs** We generate 200 random instances and observational distributions of the DTR described in Fig. 1. We assume treatments  $X_1, X_2$ , states  $S_1, S_2$  and primary outcome  $Y$  are all binary variables; values of each variable are decided by their corresponding unobserved counterfactuals  $S_{2_{x_1}}, X_{2_{x_1}}, Y_{\bar{x}_2}$  following definitions in [4, 10]. The probabilities of the joint distribution  $P(s_1, x_1, s_{2_{x_1}}, x_{2_{x_1}}, y_{\bar{x}_2})$  are drawn randomly over  $[0, 1]$ . The cumulative regrets average among all random DTRs are reported in Fig. 2a. We find that online methods (*uc-dtr*, *uc<sup>c</sup>-dtr*) dominate randomized assignments (*rand*); RL learners that leverage causal bounds (*uc<sup>c</sup>-dtr*) consistently dominates learners that do not (*uc-dtr*). Fig. 2b reports the relative improvement in total regrets of *uc<sup>c</sup>-dtr* compared to *uc-dtr* among 200 instances: *uc<sup>c</sup>-dtr* outperforms *uc-dtr* in over 80% of generated DTRs. This suggests that causal bounds derived from the observational data are beneficial in most instances.

**Cancer Treatment** We test the survival model of the two-stage clinical trial conducted by the Cancer and Leukemia Group B [18, 39]. Protocol 8923 was a double-blind, placebo controlled two-stage trial reported by [31] examining the effects of infusions of granulocyte-macrophage colony-stimulating factor (GM-CSF) after initial chemotherapy in patients with acute myelogenous leukemia (AML). Standard chemotherapy for AML could place patients at increased risk of death



due to infection or bleeding-related complications. GM-CSF administered after chemotherapy might assist patient recovery, thus reducing the number of deaths due to such complications. Patients were randomized initially to GM-CSF or placebo following standard chemotherapy. Later, patients meeting the criteria of complete remission and consenting to further participation were offered a second randomization to one of two intensification treatments.

Fig. 1a describes the DTR of this two-stage trial.  $X_1$  represents the initial GM-CSF administration and  $X_2$  represents the intensification treatment; the initial state  $S_1 = \emptyset$  and  $S_2$  indicates the complete remission after the first treatment; the primary outcome  $Y$  indicates the survival of patients at the time of recording. We generate observational samples using *age* of patients as UCs  $U$ . The cumulative regrets average among all random DTRs are reported in Fig. 2b. We find that *rand* performs worst among all strategies; *uc-dtr* finds the optimal policy with sub-linear regrets. Interestingly, *uc<sup>c</sup>-dtr* converges almost immediately, suggesting that causal bounds derived from confounded observations could significantly improve the performance of online learners.

## 5 Conclusion

In this paper, we investigated the online reinforcement learning problem for selecting the optimal DTR provided with abundant, yet imperfect observations made about the underlying environment. We first presented an online RL algorithm with near-optimal regret bounds in DTRs solely based on the knowledge about state-action domains. We further derived causal bounds about the system dynamics in DTRs from the observational data. These bounds could be incorporated in a simple, yet principled way to improve the performance of online RL learners. In today’s healthcare, for example, the growing use of mobile devices opens new opportunities in continuous monitoring of patients’ conditions and just-in-time interventions. We believe that our results constitute a significant step towards the development of a more principled and robust science of precision medicine.

## Acknowledgments

This research is supported in parts by grants from IBM Research, Adobe Research, NSF IIS-1704352, and IIS-1750807 (CAREER).

## References

- [1] P. Auer, N. Cesa-Bianchi, and P. Fischer. Finite-time analysis of the multiarmed bandit problem. *Machine learning*, 47(2-3):235–256, 2002.
- [2] P. Auer, N. Cesa-Bianchi, Y. Freund, and R. E. Schapire. The nonstochastic multiarmed bandit problem. *SIAM journal on computing*, 32(1):48–77, 2002.
- [3] K. Azizzadenesheli, A. Lazaric, and A. Anandkumar. Reinforcement learning of pomdp’s using spectral methods. In *COLT*, 2016.
- [4] A. Balke and J. Pearl. Counterfactuals and policy analysis in structural models. In *Proceedings of the Eleventh Conference on Uncertainty in Artificial Intelligence*, pages 11–18, 1995.
- [5] E. Bareinboim and J. Pearl. Causal inference and the data-fusion problem. *Proceedings of the National Academy of Sciences*, 113:7345–7352, 2016.
- [6] R. Bellman. Dynamic programming. *Science*, 153(3731):34–37, 1966.
- [7] B. Chakraborty. Dynamic treatment regimes for managing chronic health conditions: a statistical perspective. *American journal of public health*, 101(1):40–45, 2011.
- [8] B. Chakraborty and E. Moodie. *Statistical methods for dynamic treatment regimes*. Springer, 2013.
- [9] B. Chakraborty and S. A. Murphy. Dynamic treatment regimes. *Annual review of statistics and its application*, 1:447–464, 2014.
- [10] C. Frangakis and D. Rubin. Principal stratification in causal inference. *Biometrics*, 1(58):21–29, 2002.

- [11] Z. D. Guo, S. Doroudi, and E. Brunskill. A pac rl algorithm for episodic pomdps. In *Artificial Intelligence and Statistics*, pages 510–518, 2016.
- [12] W. HOEFFDING. Probability inequalities for sums of bounded random variables. *J. Am. Stat. Associ.*, 58(301):13–30, 1963.
- [13] T. Jaksch, R. Ortner, and P. Auer. Near-optimal regret bounds for reinforcement learning. *Journal of Machine Learning Research*, 11(Apr):1563–1600, 2010.
- [14] P. W. Lavori and R. Dawson. A design for testing clinical strategies: biased adaptive within-subject randomization. *Journal of the Royal Statistical Society: Series A (Statistics in Society)*, 163(1):29–38, 2000.
- [15] P. W. Lavori and R. Dawson. Adaptive treatment strategies in chronic disease. *Annu. Rev. Med.*, 59:443–453, 2008.
- [16] S. Lee, J. D. Correa, and E. Bareinboim. General identifiability with arbitrary surrogate experiments. In *Proceedings of Thirty-fifth Conference on Uncertainty in Artificial Intelligence (UAI)*, Corvallis, OR, 2019. AUAI Press.
- [17] Q. Liu and A. Ihler. Belief propagation for structured decision making. In *Proceedings of the Twenty-Eighth Conference on Uncertainty in Artificial Intelligence*, pages 523–532. AUAI Press, 2012.
- [18] J. K. Lunceford, M. Davidian, and A. A. Tsiatis. Estimation of survival distributions of treatment policies in two-stage randomization designs in clinical trials. *Biometrics*, 58(1):48–57, 2002.
- [19] C. Manski. Nonparametric bounds on treatment effects. *American Economic Review, Papers and Proceedings*, 80:319–323, 1990.
- [20] S. A. Murphy. Optimal dynamic treatment regimes. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 65(2):331–355, 2003.
- [21] S. A. Murphy. An experimental design for the development of adaptive treatment strategies. *Statistics in medicine*, 24(10):1455–1481, 2005.
- [22] S. A. Murphy. A generalization error for q-learning. *Journal of Machine Learning Research*, 6(Jul):1073–1097, 2005.
- [23] S. A. Murphy, M. J. van der Laan, J. M. Robins, and C. P. P. R. Group. Marginal mean models for dynamic regimes. *Journal of the American Statistical Association*, 96(456):1410–1423, 2001.
- [24] I. Osband and B. Van Roy. Near-optimal reinforcement learning in factored mdps. In *Advances in Neural Information Processing Systems*, pages 604–612, 2014.
- [25] J. Pearl. *Causality: Models, Reasoning, and Inference*. Cambridge University Press, New York, 2000. 2nd edition, 2009.
- [26] J. Pearl and J. Robins. Probabilistic evaluation of sequential plans from causal models with hidden variables. In P. Besnard and S. Hanks, editors, *Uncertainty in Artificial Intelligence 11*, pages 444–453. Morgan Kaufmann, San Francisco, 1995.
- [27] J. Robins, L. Orellana, and A. Rotnitzky. Estimation and extrapolation of optimal treatment and testing strategies. *Statistics in medicine*, 27(23):4678–4721, 2008.
- [28] P. Rosenbaum and D. Rubin. The central role of propensity score in observational studies for causal effects. *Biometrika*, 70:41–55, 1983.
- [29] D. Rubin. Bayesian inference for causal effects: The role of randomization. *Annals of Statistics*, 6(1):34–58, 1978.
- [30] P. Spirtes, C. N. Glymour, and R. Scheines. *Causation, prediction, and search*. MIT press, 2000.

- [31] R. M. Stone, D. T. Berg, S. L. George, R. K. Dodge, P. A. Paciucci, P. Schulman, E. J. Lee, J. O. Moore, B. L. Powell, and C. A. Schiffer. Granulocyte–macrophage colony-stimulating factor after initial chemotherapy for elderly patients with primary acute myelogenous leukemia. *New England Journal of Medicine*, 332(25):1671–1677, 1995.
- [32] A. L. Strehl, L. Li, E. Wiewiora, J. Langford, and M. L. Littman. Pac model-free reinforcement learning. In *Proceedings of the 23rd international conference on Machine learning*, pages 881–888. ACM, 2006.
- [33] R. S. Sutton and A. G. Barto. *Reinforcement learning: An introduction*. MIT press, 1998.
- [34] I. Szita and C. Szepesvári. Model-based reinforcement learning with nearly tight exploration complexity bounds. In *Proceedings of the 27th International Conference on Machine Learning (ICML-10)*, pages 1031–1038, 2010.
- [35] P. F. Thall, R. E. Millikan, and H.-G. Sung. Evaluating multiple treatment courses in clinical trials. *Statistics in medicine*, 19(8):1011–1028, 2000.
- [36] P. F. Thall, H.-G. Sung, and E. H. Estey. Selecting therapeutic strategies based on efficacy and death in multicourse clinical trials. *Journal of the American Statistical Association*, 97(457):29–39, 2002.
- [37] W. R. Thompson. On the likelihood that one unknown probability exceeds another in view of the evidence of two samples. *Biometrika*, 25(3/4):285–294, 1933.
- [38] E. H. Wagner, B. T. Austin, C. Davis, M. Hindmarsh, J. Schaefer, and A. Bonomi. Improving chronic illness care: translating evidence into action. *Health affairs*, 20(6):64–78, 2001.
- [39] A. S. Wahed and A. A. Tsiatis. Optimal estimator for the survival distribution and related quantities for treatment policies in two-stage randomization designs in clinical trials. *Biometrics*, 60(1):124–133, 2004.
- [40] A. S. Wahed and A. A. Tsiatis. Semiparametric efficient estimation of survival distributions in two-stage randomisation designs in clinical trials with censored data. *Biometrika*, 93(1):163–177, 2006.
- [41] J. Zhang and E. Bareinboim. Transfer learning in multi-armed bandits: a causal approach. In *Proceedings of the 26th International Joint Conference on Artificial Intelligence*, pages 1340–1346. AAAI Press, 2017.

---

# “Near-Optimal Reinforcement Learning in Dynamic Treatment Regimes”

## Supplemental Material

---

### Appendix I. Proofs

In this section, we provide proofs for the theoretical results presented in the main text.

#### Proof of Theorems 1 to 3

We start by introducing necessary notations for the proof. We say an episode  $t$  is  $\epsilon$ -bad if  $V_{\pi^*}(M^*) - Y^t \geq \epsilon$ . Let  $T_\epsilon$  be the number of episodes taken by UC-DTR that are  $\epsilon$ -bad. Let  $L_\epsilon$  denote the indices of the  $\epsilon$ -bad episodes up to episode  $T$ . The cumulative regret  $R_\epsilon(T)$  in  $\epsilon$ -bad episodes up to episode  $T$  is defined as  $R_\epsilon(T) = \sum_{t \in L_\epsilon} V_{\pi^*}(M^*) - Y^t$ . For any  $k = 1, \dots, K$ , we define event counts  $N(\bar{s}_k, \bar{x}_k)$  in total episodes  $T$  as  $N(\bar{s}_k, \bar{x}_k) = \sum_{t=1}^T I_{\bar{s}_k^t = \bar{s}_k, \bar{x}_k^t = \bar{x}_k}$ . Finally, we denote by  $\mathcal{H}^t$  the history up to episode  $t$ , i.e.,  $\mathcal{H}^t = \{\bar{X}_K^1, \bar{S}_K^1, Y^1, \dots, \bar{X}_K^t, \bar{S}_K^t, Y^t\}$ .

**Lemma 2.** Fix  $\delta \in (0, 1)$ , with probability at least  $1 - \delta$ ,

$$\sum_{t \in L_\epsilon} (E_{\bar{x}_K^t} [Y | \bar{S}_K^t] - Y^t) \leq \sqrt{\frac{T_\epsilon \log(1/\delta)}{2}}.$$

*Proof.* Let  $\mathbf{D}^T$  denote the sequence  $\{\bar{X}_K^1, \bar{S}_K^1, \dots, \bar{X}_K^T, \bar{S}_K^T\}$ . Rewards  $Y^t$  are independent variables by conditioning on  $\mathbf{D}^T = \mathbf{d}^T$ . Applying Hoeffding’s inequality gives:

$$P\left(\sum_{t \in L_\epsilon} (E_{\bar{x}_K^t} [Y | \bar{S}_K^t] - Y^t) \geq \sqrt{\frac{T_\epsilon \log(1/\delta)}{2}} \mid \mathbf{d}^T\right) \leq \delta.$$

We thus have:

$$P\left(\sum_{t \in L_\epsilon} (E_{\bar{x}_K^t} [Y | \bar{S}_K^t] - Y^t) \geq \sqrt{\frac{T_\epsilon \log(1/\delta)}{2}} \mid \mathbf{d}^T\right) \leq \delta \sum_{\mathbf{d}^T} P(\mathbf{d}^T) = \delta. \quad \square$$

**Lemma 3.** Fix  $\epsilon > 0$ ,  $\delta \in (0, 1)$ . With probability (w.p.) of at least  $1 - \delta$ , it holds for any  $T > 1$ ,  $R_\epsilon(T)$  of UC-DTR with parameter  $\delta$  is bounded by

$$R_\epsilon(T) \leq 12K\sqrt{|\mathcal{S}||\mathcal{X}|T_\epsilon \log(2K|\mathcal{S}||\mathcal{X}|T/\delta)} + 4K\sqrt{T_\epsilon \log(2T/\delta)}$$

*Proof.* Let  $M^*$  denote the underlying DTR. Recall that  $\mathcal{M}_t$  is a set of DTR instances such that for any  $M \in \mathcal{M}_t$ , its system dynamics satisfy

$$\left\| P_{\bar{x}_k}^M(\cdot | \bar{s}_k) - \hat{P}_{\bar{x}_k}^t(\cdot | \bar{s}_k) \right\|_1 \leq \sqrt{\frac{6|\mathcal{S}_{k+1}| \log(2K|\mathcal{S}_k||\mathcal{X}_k|t/\delta)}{\max\{1, N^t(\bar{s}_k, \bar{x}_k)\}}}, \quad (17)$$

$$\left| E_{\bar{x}_K}^M [Y | \bar{s}_K] - \hat{E}_{\bar{x}_K}^t [Y | \bar{s}_K] \right| \leq \sqrt{\frac{2 \log(2K|\mathcal{S}||\mathcal{X}|t/\delta)}{\max\{1, N^t(\bar{s}_K, \bar{x}_K)\}}}. \quad (18)$$

By union bounds and Hoeffding's inequality (following a similar argument in [4, C.1]),

$$P(M^* \notin \mathcal{M}_t) \leq \frac{\delta}{4t^2}.$$

Since  $\sum_{t=1}^{\infty} \frac{1}{4t^2} \leq \frac{\pi^2}{24} \delta < \frac{\delta}{2}$ , it follows that with probability at least  $1 - \frac{\delta}{2}$ ,  $M^* \in \mathcal{M}_t$  for all episodes  $t = 1, 2, \dots$ .

For the remainder of the proof, we will assume that  $M^* \in \mathcal{M}_t$  for all  $t$ . Let  $E_{\bar{\mathbf{x}}_K}^{M_t}[Y|\bar{\mathbf{s}}_K]$  denote the conditional expected reward in the optimistic DTR  $M_t$ . We can write  $R_\epsilon(T)$  as:

$$R_\epsilon(T) = \sum_{t \in L_\epsilon} (V_{\pi^*}(M^*) - E_{\bar{\mathbf{x}}_K}^{M_t}[Y|\bar{\mathbf{s}}_K^t]) \quad (19)$$

$$+ \sum_{t \in L_\epsilon} (E_{\bar{\mathbf{x}}_K}^{M_t}[Y|\bar{\mathbf{s}}_K^t] - E_{\bar{\mathbf{x}}_K}[\bar{Y}|\bar{\mathbf{s}}_K^t]) \quad (20)$$

$$+ \sum_{t \in L_\epsilon} (E_{\bar{\mathbf{x}}_K}[\bar{Y}|\bar{\mathbf{s}}_K^t] - Y^t). \quad (21)$$

We will next derive bounds over  $R_\epsilon(T)$  by bounding quantities in Eqs. (19) to (21) separately.

**Bounding Eq. (19)** For any DTR  $M$  and policy  $\pi$ , let  $V_\pi(\bar{\mathbf{s}}_k, \bar{\mathbf{x}}_{k-1}; M) = E_\pi^M[Y|\bar{\mathbf{s}}_k, \bar{\mathbf{x}}_{k-1}]$  and  $V_\pi(\bar{\mathbf{s}}_k, \bar{\mathbf{x}}_k; M) = E_\pi^M[Y|\bar{\mathbf{s}}_k, \bar{\mathbf{x}}_k]$ . Since  $M^* \in \mathcal{M}_t$ , we must have  $V_{\pi^*}(s_1; M^*) \leq V_{\pi_t}(s_1; M_t)$ , i.e., the maximal expected reward of the optimal reward in the optimistic  $M_t$  is no less than that in the underlying DTR  $M^*$  for any initial state  $s_1$ . Further, since  $\pi_t$  is deterministic, for any stage  $k$  and DTR  $M$ ,

$$V_{\pi_t}(\bar{\mathbf{s}}_k^t, \bar{\mathbf{x}}_{k-1}^t; M) = V_{\pi_t}(\bar{\mathbf{s}}_k^t, \bar{\mathbf{x}}_k^t; M). \quad (22)$$

We thus have

$$V_{\pi^*}(M^*) - E_{\bar{\mathbf{x}}_K}^{M_t}[Y|\bar{\mathbf{s}}_K^t] \leq V_{\pi^*}(M^*) - V_{\pi^*}(\bar{\mathbf{s}}_1^t; M^*) + V_{\pi_t}(\bar{\mathbf{s}}_1^t, \bar{\mathbf{x}}_1^t; M^*) - E_{\bar{\mathbf{x}}_K}^{M_t}[Y|\bar{\mathbf{s}}_K^t].$$

Let  $M_t(k)$  denote a combined DTR obtained from  $M^*$  and  $M_t$  such that

- for  $i = 0, 1, \dots, k-1$ , its transition probability  $P_{\bar{\mathbf{x}}_i}^{M_t(k)}(s_{i+1}|\bar{\mathbf{s}}_i)$  coincides with the transition probability  $P_{\bar{\mathbf{x}}_i}(s_{i+1}|\bar{\mathbf{s}}_i)$  in the real DTR  $M^*$ ;
- for  $i = k, \dots, K-1$ , its transition probability  $P_{\bar{\mathbf{x}}_i}^{M_t(k)}(s_{i+1}|\bar{\mathbf{s}}_i)$  coincides with the transition probability  $P_{\bar{\mathbf{x}}_i}^{M_t}(s_{i+1}|\bar{\mathbf{s}}_i)$  in the optimistic  $M_t$

This is, for any  $\pi \in \Pi$ , the interventional distribution  $P_\pi^{M_t(k)}(\bar{\mathbf{x}}_K, \bar{\mathbf{s}}_K, y)$  factorizes as follows:

$$\begin{aligned} P_\pi^{M_t(k)}(\bar{\mathbf{x}}_K, \bar{\mathbf{s}}_K, y) &= P_{\bar{\mathbf{x}}_K}^{M_t}(y|\bar{\mathbf{s}}_K) \prod_{i=0}^{k-1} P_{\bar{\mathbf{x}}_i}(s_{i+1}|\bar{\mathbf{s}}_i) \\ &\quad \cdot \prod_{j=k}^{K-1} P_{\bar{\mathbf{x}}_j}^{M_t}(s_{j+1}|\bar{\mathbf{s}}_j) \prod_{l=1}^{K-1} \pi_{l+1}(x_{l+1}|\bar{\mathbf{s}}_{l+1}, \bar{\mathbf{x}}_l). \end{aligned} \quad (23)$$

Obviously,  $E_{\bar{\mathbf{x}}_K}^{M_t}[Y|\bar{\mathbf{s}}_K^t] = V_{\pi_t}(\bar{\mathbf{s}}_K^t, \bar{\mathbf{x}}_K^t; M_t^{(K)})$  and  $V_{\pi_t}(\bar{\mathbf{s}}_1^t, \bar{\mathbf{x}}_1^t; M_t) = V_{\pi_t}(s_1^t, x_1^t; M_t^{(1)})$ . We thus have

$$\begin{aligned} V_{\pi_t}(\bar{\mathbf{s}}_1^t, \bar{\mathbf{x}}_1^t; M_t) - E_{\bar{\mathbf{x}}_K}^{M_t}[Y|\bar{\mathbf{s}}_K^t] &= V_{\pi_t}(\bar{\mathbf{s}}_1^t, \bar{\mathbf{x}}_1^t; M_t^{(1)}) - V_{\pi_t}(\bar{\mathbf{s}}_K^t, \bar{\mathbf{x}}_K^t; M_t^{(K)}) \\ &= \sum_{k=1}^{K-1} V_{\pi_t}(\bar{\mathbf{s}}_k^t, \bar{\mathbf{x}}_k^t; M_t^{(1)}) - V_{\pi_t}(\bar{\mathbf{s}}_{k+1}^t, \bar{\mathbf{x}}_{k+1}^t; M_t^{(K)}) \\ &= \sum_{k=1}^{K-1} V_{\pi_t}(\bar{\mathbf{s}}_k^t, \bar{\mathbf{x}}_k^t; M_t^{(1)}) - V_{\pi_t}(\bar{\mathbf{s}}_{k+1}^t, \bar{\mathbf{x}}_k^t; M_t^{(K)}). \end{aligned}$$

The last step is ensured by Eq. (22). We further have:

$$\begin{aligned} V_{\pi_t}(\bar{\mathbf{S}}_1^t, \bar{\mathbf{X}}_1^t; M_t) - E_{\bar{\mathbf{X}}_K^t}^{M_t}[Y|\bar{\mathbf{S}}_K^t] &= \sum_{k=1}^{K-1} V_{\pi_t}(\bar{\mathbf{S}}_k^t, \bar{\mathbf{X}}_k^t; M_t^{(k)}) - V_{\pi_t}(\bar{\mathbf{S}}_k^t, \bar{\mathbf{X}}_k^t; M_t^{(k+1)}) \\ &\quad + \sum_{k=1}^{K-1} V_{\pi_t}(\bar{\mathbf{S}}_k^t, \bar{\mathbf{X}}_k^t; M_t^{(k+1)}) - V_{\pi_t}(\bar{\mathbf{S}}_{k+1}^t, \bar{\mathbf{X}}_k^t; M_t^{(k+1)}). \end{aligned}$$

Eq. (19) can thus be written as:

$$\sum_{t \in L_\epsilon} (V_{\pi_t}(M_t) - E_{\bar{\mathbf{X}}_K^t}^{M_t}[Y|\bar{\mathbf{S}}_K^t]) = \sum_{k=1}^{K-1} \sum_{t \in L_\epsilon} V_{\pi_t}(\bar{\mathbf{S}}_k^t, \bar{\mathbf{X}}_k^t; M_t^{(k)}) - V_{\pi_t}(\bar{\mathbf{S}}_k^t, \bar{\mathbf{X}}_k^t; M_t^{(k+1)}) + \sum_{t \in L_\epsilon} Z_t,$$

where  $Z_t$  is defined as

$$Z_t = V_{\pi^*}(M^*) - V_{\pi^*}(\bar{\mathbf{S}}_1^t; M) + \sum_{k=1}^{K-1} V_{\pi_t}(\bar{\mathbf{S}}_k^t, \bar{\mathbf{X}}_k^t; M_t^{(k+1)}) - V_{\pi_t}(\bar{\mathbf{S}}_{k+1}^t, \bar{\mathbf{X}}_k^t; M_t^{(k+1)})$$

By Eq. (23) and basic probabilistic operations,

$$\begin{aligned} &V_{\pi_t}(\bar{\mathbf{S}}_k^t, \bar{\mathbf{X}}_k^t; M_t^{(k)}) - V_{\pi_t}(\bar{\mathbf{S}}_k^t, \bar{\mathbf{X}}_k^t; M_t^{(k+1)}) \\ &= \sum_{s_{k+1}} (P^{M_t}(s_{k+1}|\bar{\mathbf{S}}_k, \bar{\mathbf{X}}_k) - P(s_{k+1}|\bar{\mathbf{S}}_k, \bar{\mathbf{X}}_k)) V_{\pi_t}(s_{k+1}, \bar{\mathbf{S}}_k^t, \bar{\mathbf{X}}_k^t; M_t) \\ &\leq \left\| P_{\bar{\mathbf{x}}_k}^{M_t}(\cdot|\bar{\mathbf{s}}_k) - P_{\bar{\mathbf{x}}_k}(\cdot|\bar{\mathbf{s}}_k) \right\|_1 \max_{s_{k+1}} V_{\pi_t}(s_{k+1}, \bar{\mathbf{S}}_k^t, \bar{\mathbf{X}}_k^t; M_t) \\ &\leq 2\sqrt{6|\mathcal{S}_{k+1}| \log(2K|\bar{\mathbf{S}}_k||\bar{\mathbf{X}}_k|T/\delta)} \frac{1}{\sqrt{\max\{1, N^t(\bar{\mathbf{S}}_k^t, \bar{\mathbf{X}}_k^t)\}}} \end{aligned}$$

The last step follows from Eq. (17). From results in [4, D], we have

$$\sum_{t \in L_\epsilon} \frac{1}{\sqrt{\max\{1, N^t(\bar{\mathbf{S}}_k^t, \bar{\mathbf{X}}_k^t)\}}} \leq (\sqrt{2} + 1) \sqrt{T_\epsilon |\bar{\mathbf{S}}_k| |\bar{\mathbf{X}}_k|}.$$

This implies:

$$\begin{aligned} &\sum_{t \in L_\epsilon} \sum_{k=1}^{K-1} V_{\pi_t}(\bar{\mathbf{S}}_k^t, \bar{\mathbf{X}}_k^t; M_t^{(k)}) - V_{\pi_t}(\bar{\mathbf{S}}_k^t, \bar{\mathbf{X}}_k^t; M_t^{(k+1)}) \\ &\leq \sum_{k=1}^{K-1} 2(\sqrt{2} + 1) \sqrt{6T_\epsilon |\bar{\mathbf{S}}_{k+1}| |\bar{\mathbf{X}}_k| \log(2K|\bar{\mathbf{S}}_k||\bar{\mathbf{X}}_k|T/\delta)} \\ &\leq 2(\sqrt{2} + 1)(K-1) \sqrt{6T_\epsilon |\mathcal{S}| |\mathcal{X}| \log(2K|\mathcal{S}| |\mathcal{X}| T/\delta)} \end{aligned} \quad (24)$$

Let  $\mathcal{H}^t$  denote the history up to episode  $t$ , i.e.,  $\{\bar{\mathbf{X}}_K^1, \bar{\mathbf{S}}_K^1, Y^1, \dots, \bar{\mathbf{X}}_K^t, \bar{\mathbf{S}}_K^t, Y^t\}$ . Since  $|Z_t| \leq K$  and  $E[Z_{t+1}|\mathcal{H}_t] = 0$ ,  $\{Z_t : t \in L_\epsilon\}$  is a sequence of martingale differences. By Azuma-Hoeffding inequality [3], we have, with probability at least  $1 - \frac{\delta}{8T^2}$ ,

$$\sum_{t \in L_\epsilon} Z_t \leq K \sqrt{6T_\epsilon \log(2T/\delta)} \quad (25)$$

Since  $\sum_{T=1}^\infty \frac{1}{8T^2} \leq \frac{\pi^2}{48} \delta < \frac{\delta}{4}$ , the above inequality holds with probability  $1 - \frac{\delta}{4}$  for all  $T > 1$ . Eqs. (24) and (25) combined give

$$\begin{aligned} &\sum_{t \in L_\epsilon} (V_{\pi^*}(M^*) - E_{\bar{\mathbf{X}}_K^t}^{M_t}[Y|\bar{\mathbf{S}}_K^t]) \\ &\leq 2(\sqrt{2} + 1)(K-1) \sqrt{6T_\epsilon |\mathcal{S}| |\mathcal{X}| \log(2K|\mathcal{S}| |\mathcal{X}| T/\delta)} + K \sqrt{6T_\epsilon \log(2T/\delta)} \end{aligned} \quad (26)$$

**Bounding Eq. (20)** Since both  $M^*, M_t$  are in the set  $\mathcal{M}_t$ ,

$$\begin{aligned} E_{\bar{\mathbf{X}}_K^t}^{M_t}[Y|\bar{\mathbf{S}}_K^t] - E_{\bar{\mathbf{X}}_K^t}[Y|\bar{\mathbf{S}}_K^t] &\leq \left| E_{\bar{\mathbf{x}}_K}^{M_t}[Y|\bar{\mathbf{s}}_K] - \hat{E}_{\bar{\mathbf{x}}_K}^t[Y|\bar{\mathbf{s}}_K] \right| + \left| E_{\bar{\mathbf{X}}_K^t}[Y|\bar{\mathbf{S}}_K^t] - \hat{E}_{\bar{\mathbf{x}}_K}^t[Y|\bar{\mathbf{s}}_K] \right| \\ &\leq 2\sqrt{2\log(2K|\mathcal{S}||\mathcal{X}|T/\delta)} \frac{1}{\sqrt{\max\{1, N^t(\bar{\mathbf{S}}_K^t, \bar{\mathbf{X}}_K^t)\}}} \end{aligned}$$

The last step follows from Eq. (18). From results in [4, D], we have

$$\sum_{t \in L_\epsilon} \frac{1}{\sqrt{\max\{1, N^t(\bar{\mathbf{S}}_K^t, \bar{\mathbf{X}}_K^t)\}}} \leq (\sqrt{2} + 1)\sqrt{T_\epsilon|\mathcal{S}||\mathcal{X}|}.$$

This implies

$$\sum_{t \in L_\epsilon} (E_{\bar{\mathbf{X}}_K^t}^{M_t}[Y|\bar{\mathbf{S}}_K^t] - E_{\bar{\mathbf{X}}_K^t}[Y|\bar{\mathbf{S}}_K^t]) \leq 2(\sqrt{2} + 1)\sqrt{2T_\epsilon|\mathcal{S}||\mathcal{X}|\log(2K|\mathcal{S}||\mathcal{X}|T/\delta)} \quad (27)$$

**Bounding Eq. (21)** By Lem. 2, we have with probability at least  $1 - \frac{\delta}{8T^2}$ ,

$$\sum_{t \in L_\epsilon} (E_{\bar{\mathbf{X}}_K^t}[Y|\bar{\mathbf{S}}_K^t] - Y^t) \leq \sqrt{\frac{3T_\epsilon \log(2T/\delta)}{2}} \quad (28)$$

Since  $\sum_{T=1}^\infty \frac{1}{8T^2} \leq \frac{\pi^2}{48} \delta < \frac{\delta}{4}$ , the above equation holds with probability  $1 - \frac{\delta}{4}$  for any  $T$ .

Eqs. (26) to (28) together give that, with probability at least  $1 - \frac{\delta}{2} - \frac{\delta}{4} - \frac{\delta}{4} = 1 - \delta$ ,

$$\begin{aligned} R_\epsilon(T) &\leq (K-1)2(\sqrt{2} + 1)\sqrt{6T_\epsilon|\mathcal{S}||\mathcal{X}|\log(2K|\mathcal{S}||\mathcal{X}|T/\delta)} + K\sqrt{6T_\epsilon \log(2T/\delta)} \\ &\quad + 2(\sqrt{2} + 1)\sqrt{2T_\epsilon|\mathcal{S}||\mathcal{X}|\log(2K|\mathcal{S}||\mathcal{X}|T/\delta)} + \sqrt{\frac{3T_\epsilon \log(2T/\delta)}{2}}. \end{aligned}$$

A quick simplification gives:

$$R_\epsilon(T) \leq 12K\sqrt{|\mathcal{S}||\mathcal{X}|T_\epsilon \log(2K|\mathcal{S}||\mathcal{X}|T/\delta)} + 4K\sqrt{T_\epsilon \log(2T/\delta)}. \quad \square$$

**Theorem 1.** Fix a  $\delta \in (0, 1)$ . With probability (w.p.) of at least  $1 - \delta$ , it holds for any  $T > 1$ , the regret of UC-DTR with parameter  $\delta$  is bounded by

$$R(T) \leq 12K\sqrt{|\mathcal{S}||\mathcal{X}|T \log(2K|\mathcal{S}||\mathcal{X}|T/\delta)} + 4K\sqrt{T \log(2T/\delta)}.$$

*Proof.* Fix  $\epsilon = 0$ . Naturally,  $T_\epsilon = T$  and  $R_\epsilon(T) = R(T)$ . By Lem. 3,

$$R(T) \leq 12K\sqrt{|\mathcal{S}||\mathcal{X}|T \log(2K|\mathcal{S}||\mathcal{X}|T/\delta)} + 4K\sqrt{T \log(2T/\delta)}. \quad \square$$

**Theorem 2.** For any  $T \geq 1$ , with parameter  $\delta = \frac{1}{T}$ , the expected regret of UC-DTR is bounded by

$$E[R(T)] \leq \max_{\pi \in \Pi^-} \left\{ \frac{33^2 K^2 |\mathcal{S}||\mathcal{X}| \log(T)}{\Delta_\pi} + \frac{32}{\Delta_\pi^3} + \frac{4}{\Delta_\pi} \right\} + 1.$$

*Proof.* By Lem. 3 and a quick simplification, we have

$$R_\epsilon(T) \leq 23K\sqrt{|\mathcal{S}||\mathcal{X}|T_\epsilon \log(T/\delta)}.$$

Since  $R_\epsilon(T) \geq \epsilon T_\epsilon$ ,  $\epsilon T_\epsilon \leq 23K\sqrt{|\mathcal{S}||\mathcal{X}|T_\epsilon \log(T/\delta)}$ , which implies

$$T_\epsilon \leq \frac{23^2 K^2 |\mathcal{S}||\mathcal{X}| \log(T/\delta)}{\epsilon^2}. \quad (29)$$

This implies that, with probability at least  $1 - \delta$ ,

$$R_\epsilon(T) \leq 23K\sqrt{|\mathcal{S}||\mathcal{X}|T_\epsilon \log(T/\delta)} = \frac{23^2 K^2 |\mathcal{S}||\mathcal{X}| \log(T/\delta)}{\epsilon}$$

Let  $\Delta = \arg \min_{\pi \in \Pi} \Delta_{\pi}$ . Fix  $\epsilon = \frac{\Delta}{2}$ ,  $\delta = \frac{1}{T}$ , we have

$$E[R_{\frac{\Delta}{2}}(T)] \leq \frac{33^2 K^2 |\mathcal{S}| |\mathcal{X}| \log(T)}{\Delta} + 1. \quad (30)$$

We now only need to bound the regrets cumulated in the episodes that are not  $\epsilon$ -bad, which we call  $\epsilon$ -good. Let  $\tilde{R}_{\epsilon}(T)$  denote the regret in episodes that are  $\epsilon$ -good. Let  $\tilde{T}_{\epsilon}$  denote the total number of  $\epsilon$ -good episodes and let  $\tilde{L}_{\epsilon}$  be indices of  $\epsilon$ -good episodes. Fix  $\epsilon = \frac{\Delta}{2}$ , for any  $\epsilon$ -good episode  $t$ , we have  $V_{\pi_t}(M^*) - Y^t < \epsilon$ . Fix event  $\tilde{T}_{\frac{\Delta}{2}} = t$ ,

$$\tilde{R}_{\epsilon}(T) = \sum_{i \in \tilde{L}_{\epsilon}} V_{\pi^*}(M^*) - Y^i \leq t \frac{\Delta}{2}.$$

The above inequality is equivalent to

$$\begin{aligned} \sum_{i \in \tilde{L}_{\frac{\Delta}{2}}} V_{\pi^*}(M^*) - V_{\pi_i}(M^*) - Y^i &\leq t \frac{\Delta}{2} - \sum_{i \in \tilde{L}_{\frac{\Delta}{2}}} V_{\pi_i}(M^*) \\ \Rightarrow \sum_{i \in \tilde{L}_{\frac{\Delta}{2}}} \Delta_{\pi_i} - Y^i &\leq t \frac{\Delta}{2} - \sum_{i \in \tilde{L}_{\frac{\Delta}{2}}} V_{\pi_i}(M^*) \\ \Rightarrow \sum_{i \in \tilde{L}_{\frac{\Delta}{2}}} \Delta - Y^i &\leq t \frac{\Delta}{2} - \sum_{i \in \tilde{L}_{\frac{\Delta}{2}}} V_{\pi_i}(M^*) \end{aligned}$$

Since  $|\tilde{L}_{\epsilon}| = \tilde{T}_{\frac{\Delta}{2}}$ , we have

$$\tilde{T}_{\frac{\Delta}{2}} = t \Rightarrow \sum_{i \in \tilde{L}_{\frac{\Delta}{2}}} V_{\pi_i}(M^*) - Y^i \leq -t \frac{\Delta}{2}. \quad (31)$$

We could thus bound  $E[\tilde{R}_{\frac{\Delta}{2}}(T)]$  as

$$E[\tilde{R}_{\frac{\Delta}{2}}(T)] \leq \frac{\Delta}{2} E[\tilde{T}_{\frac{\Delta}{2}}(T)] \leq \frac{\Delta}{2} \sum_{t=1}^T t P(\tilde{T}_{\frac{\Delta}{2}} = t)$$

By Eq. (31), we further have

$$E[\tilde{R}_{\frac{\Delta}{2}}(T)] \leq \frac{\Delta}{2} \sum_{t=1}^T t P\left(\sum_{i \in \tilde{L}_{\frac{\Delta}{2}}} V_{\pi_i}(M^*) - Y^i \leq -t \frac{\Delta}{2}\right)$$

Let  $C_t = V_{\pi_t}(M^*) - Y^t$ . Since  $|C_t| < 1$  and  $E[C_{t+1} | \mathcal{H}^t] = 0$ ,  $\{C_i : i \in \tilde{L}_{\frac{\Delta}{2}}\}$  is a sequence of martingale differences. Applying Azuma-Hoeffding lemma gives,

$$P\left(\sum_{i \in \tilde{L}_{\frac{\Delta}{2}}} C_i \leq -t \frac{\Delta}{2}\right) \leq e^{-\frac{\Delta^2 t}{8}}.$$

Thus

$$E[\tilde{R}_{\frac{\Delta}{2}}(T)] \leq \frac{\Delta}{2} \sum_{t=1}^T t e^{-\frac{\Delta^2 t}{8}} \leq \frac{\Delta}{2} \frac{64}{\Delta^4} \left(\frac{\Delta^2}{8} + 1\right) e^{-\frac{\Delta^2}{8}}$$

which implies

$$E[\tilde{R}_{\frac{\Delta}{2}}(T)] \leq \frac{32}{\Delta^3} + \frac{4}{\Delta}. \quad (32)$$

Eqs. (30) and (32) together give:

$$E[R(T)] = E[R_{\frac{\Delta}{2}}(T)] + E[\tilde{R}_{\frac{\Delta}{2}}(T)] \leq \frac{33^2 K^2 |\mathcal{S}| |\mathcal{X}| \log(T)}{\Delta} + \frac{32}{\Delta^3} + \frac{4}{\Delta} + 1$$

The right-hand side of the above inequality is a decreasing function regarding the gap  $\Delta$ . By a quick simplification, we prove the statement.  $\square$



**Theorem 3.** For any algorithm  $\mathcal{A}$ , any natural numbers  $K \geq 1$ , and  $|\mathcal{S}^k| \geq 2, |\mathcal{X}^k| \geq 2$  for any  $k \in \{1, \dots, K\}$ , there is a DTR  $M$  with horizon  $K$ , state domains  $\mathcal{S}$  and action domains  $\mathcal{X}$ , such that the expected regret of  $\mathcal{A}$  after  $T \geq |\mathcal{S}||\mathcal{X}|$  episodes is at least

$$E[R(T)] \geq 0.05 \sqrt{|\mathcal{S}||\mathcal{X}|T}.$$

*Proof.* The classic results in bandit literature [1, Thm. 5.1] shows that for each state sequence  $K$ , there exists a bandit instance such that for any the total regret of any algorithm is lower bound by

$$E[R(T)] \geq 0.05 \sum_{\bar{s}_K} \sqrt{N(\bar{s}_K)|\mathcal{X}|},$$

where  $N(\bar{s}_K)$  is the event count  $\bar{S}_K = \bar{s}_K$  for all  $T$  episodes. The lower bound in Thm. 3 is achieved when all states  $K$  are decided uniformly at random, i.e.,  $N(\bar{s}_K) = T/|\bar{S}_K|$ .  $\square$

### Proofs of Theorems 4 to 6, Lemma 1, and Corollary 2

In this section, we provide proofs for the bounds on transition probabilities of DTRs. Our proofs build on the notion of counterfactual variables [6, Ch. 7.1] and axioms of “composition, effectiveness and reversibility” defined in [6, Ch. 7.3.1].

For a SCM  $M$ , arbitrary subsets of endogenous variables  $\mathbf{X}, \mathbf{Y}$ , the potential outcome of  $\mathbf{Y}$  to intervention  $do(\mathbf{x})$ , denoted by  $\mathbf{Y}_x(\mathbf{u})$ , is the solution for  $\mathbf{Y}$  with  $\mathbf{U} = \mathbf{u}$  in the sub-model  $M_x$ . It can be read as the counterfactual sentence “the value that  $\mathbf{Y}$  would have obtained in situation  $\mathbf{U} = \mathbf{u}$ , had  $\mathbf{X}$  been  $\mathbf{x}$ .” Statistically, averaging  $\mathbf{u}$  over the distribution  $P(\mathbf{u})$  leads to the counterfactual variables  $\mathbf{Y}_x$ . We denote  $P(\mathbf{Y}_x)$  a distribution over counterfactual variables  $\mathbf{Y}_x$ . We use  $P(\mathbf{y}_x)$  as a shorthand for probabilities  $P(\mathbf{Y}_x = \mathbf{y})$  when the identify of the counterfactual variables is clear.

We now introduce a family of DTRs which represent the exogenous variables  $\mathbf{U}$  using partitions defined by the corresponding counterfactual variables. For any  $k = 1, \dots, K-1$ , let  $S_{k+1, \bar{\mathbf{x}}_k}$  denote a set of counterfactual variables  $\{S_{k+1, \bar{\mathbf{x}}_k} : \bar{\mathbf{x}}_k \in \bar{\mathcal{X}}_k\}$ . Similarly, let  $Y_{\bar{\mathcal{X}}_K}$  denote a set  $\{Y_{\bar{\mathbf{x}}_K} : \bar{\mathbf{x}}_K \in \bar{\mathcal{X}}_K\}$ . Further, we define  $\bar{S}_{k+1, \bar{\mathbf{x}}_k}$  a set  $\{S_1, S_{2, \bar{\mathbf{x}}_1}, \dots, S_{k+1, \bar{\mathbf{x}}_k}\}$ .

**Definition 1** (Counterfactual DTR). A counterfactual dynamic treatment regime is a DTR  $\langle \mathbf{U}, \{\bar{\mathcal{X}}_K, \bar{S}_K, Y\}, \mathbf{F}, P(\mathbf{u}) \rangle$  where for  $k = 2, \dots, K$ ,

- The exogenous variables  $\mathbf{U} = \{\bar{\mathcal{X}}_K, \bar{S}_{K, \bar{\mathcal{X}}_{K-1}}, Y_{\bar{\mathcal{X}}_K}\}$ ;
- Values of  $S_1, \bar{\mathcal{X}}_K$  are drawn from  $P(\bar{\mathcal{X}}_K, \bar{S}_{K, \bar{\mathcal{X}}_{K-1}}, Y_{\bar{\mathcal{X}}_K})$ ;
- Values of  $S_k$  are decided by a function  $S_k \leftarrow \tau_k(S_{k, \bar{\mathcal{X}}_{k-1}}, \bar{\mathcal{X}}_{k-1}) = S_{k, \bar{\mathcal{X}}_{k-1}}$ ;
- Values of  $Y$  are decided by a function  $Y \leftarrow r(Y_{\bar{\mathcal{X}}_K}, \bar{\mathcal{X}}_K) = Y_{\bar{\mathcal{X}}_K}$ .

Give observational distribution  $P(\bar{s}_K, \bar{\mathbf{x}}_K, y) > 0$ , we next construct a family of counterfactual DTRs  $\mathcal{M}_{\text{OBS}}$  that are compatible with the observational distribution, i.e., for any  $M \in \mathcal{M}_{\text{OBS}}$ ,  $P^M(\bar{s}_K, \bar{\mathbf{x}}_K, y) = P(\bar{s}_K, \bar{\mathbf{x}}_K, y)$ . First, any  $M \in \mathcal{M}_{\text{OBS}}$ , its exogenous distribution  $P^M(\bar{\mathcal{X}}_K, \bar{S}_{K, \bar{\mathcal{X}}_{K-1}}, Y_{\bar{\mathcal{X}}_K})$  must satisfy the following decomposition:

$$\begin{aligned} P^M(\bar{\mathcal{X}}_K, \bar{S}_{K, \bar{\mathcal{X}}_{K-1}}, Y_{\bar{\mathcal{X}}_K}) &= P^M(s_1) \prod_{\bar{\mathbf{x}}_K^y \in \bar{\mathcal{X}}_K} P^M(Y_{\bar{\mathbf{x}}_K^y} | \bar{S}_{K, \bar{\mathcal{X}}_{K-1}}, \bar{\mathcal{X}}_K) P^M(\bar{\mathcal{X}}_K | \bar{S}_{K, \bar{\mathcal{X}}_{K-1}}, \bar{\mathcal{X}}_{K-1}) \\ &\quad \cdot \prod_{k=1}^{K-1} \prod_{\bar{\mathbf{x}}_k^{k+1} \in \bar{\mathcal{X}}_k} P^M(S_{k+1, \bar{\mathbf{x}}_k^{k+1}} | \bar{S}_{k, \bar{\mathcal{X}}_{k-1}}, \bar{\mathbf{x}}_k) P^M(\bar{\mathcal{X}}_k | \bar{S}_{k, \bar{\mathcal{X}}_{k-1}}, \bar{\mathcal{X}}_{k-1}). \end{aligned}$$

Among quantities in the above equation, we define factors  $P^M(s_1)$  as the observational probabilities  $P(s_1)$ , i.e,  $P^M(s_1) = P(s_1)$ . We further define conditional probabilities

$$\begin{aligned} P^M(y_{\bar{\mathbf{x}}_K} | \bar{S}_{K, \bar{\mathcal{X}}_{K-1}}, \bar{\mathbf{x}}_K) &= P(y | \bar{s}_K, \bar{\mathbf{x}}_K), & P^M(\bar{\mathbf{x}}_K | \bar{S}_{K, \bar{\mathcal{X}}_{K-1}}, \bar{\mathbf{x}}_{K-1}) &= P(\bar{\mathbf{x}}_K | \bar{s}_K, \bar{\mathbf{x}}_{K-1}), \\ P^M(S_{k+1, \bar{\mathbf{x}}_k} | \bar{S}_{k, \bar{\mathcal{X}}_{k-1}}, \bar{\mathbf{x}}_k) &= P(S_{k+1} | \bar{s}_k, \bar{\mathbf{x}}_k), & P^M(\bar{\mathbf{x}}_k | \bar{S}_{k, \bar{\mathcal{X}}_{k-1}}, \bar{\mathbf{x}}_{k-1}) &= P(\bar{\mathbf{x}}_k | \bar{s}_k, \bar{\mathbf{x}}_{k-1}). \end{aligned}$$

Other factors can be arbitrary conditional probabilities. It is verifiable that for any  $M \in \mathcal{M}_{\text{OBS}}$ ,  $P^M(\bar{s}_K, \bar{x}_K, y) = P(\bar{s}_K, \bar{x}_K, y)$ . To witness,

$$\begin{aligned}
P^M(\bar{s}_K, \bar{x}_K, Y) &= \sum_{k=1}^{K-1} \sum_{\{Y_{\bar{x}_K^y} : \bar{x}_K^y \neq \bar{x}_K\}} \sum_{\{S_{k+1}^{\bar{x}_K^y} : \bar{x}_K^{k+1} \neq \bar{x}_K\}} P^M(\bar{X}_K, \bar{S}_{K\bar{x}_{K-1}}, Y_{\bar{x}_K}) \\
&= P^M(s_1) \prod_{\bar{x}_K^y \in \bar{\mathcal{X}}_K} \sum_{\{Y_{\bar{x}_K^y} : \bar{x}_K^y \neq \bar{x}_K\}} P^M(Y_{\bar{x}_K^y} | \bar{S}_{K\bar{x}_{K-1}}, \bar{X}_K) P^M(\bar{X}_K | \bar{S}_{K\bar{x}_{K-1}}, \bar{x}_{K-1}) \\
&\quad \cdot \prod_{k=1}^{K-1} \prod_{\bar{x}_K^{k+1} \in \bar{\mathcal{X}}_K} \sum_{\{S_{k+1}^{\bar{x}_K^{k+1}} : \bar{x}_K^{k+1} \neq \bar{x}_K\}} P^M(S_{k+1}^{\bar{x}_K^{k+1}} | \bar{S}_{K\bar{x}_{K-1}}, \bar{X}_K) P^M(\bar{X}_K | \bar{S}_{K\bar{x}_{K-1}}, \bar{x}_{K-1}) \\
&= P^M(s_1) P^M(Y_{\bar{x}_K} | \bar{S}_{K\bar{x}_{K-1}}, \bar{X}_K) P^M(\bar{X}_K | \bar{S}_{K\bar{x}_{K-1}}, \bar{x}_{K-1}) \\
&\quad \cdot \prod_{k=1}^{K-1} P^M(S_{k+1}^{\bar{x}_K} | \bar{S}_{K\bar{x}_{K-1}}, \bar{x}_k) P^M(\bar{X}_K | \bar{S}_{K\bar{x}_{K-1}}, \bar{x}_{K-1}).
\end{aligned}$$

By definitions of  $\mathcal{M}_{\text{OBS}}$ , we thus have that, for any  $\bar{s}_K, \bar{x}_K, y$ ,

$$\begin{aligned}
P^M(\bar{s}_K, \bar{x}_K, y) &= P(s_1) P(y | \bar{s}_K, \bar{x}_K) P(\bar{x}_K | \bar{s}_K, \bar{x}_{K-1}) \prod_{k=1}^{K-1} P(s_{k+1} | \bar{s}_k, \bar{x}_k) P(\bar{x}_k | \bar{s}_k, \bar{x}_{k-1}) \\
&= P(\bar{s}_K, \bar{x}_K, y).
\end{aligned}$$

We will now use the constructions of  $\mathcal{M}_{\text{OBS}}$  to prove the non-identifiability of  $P_{\bar{x}_K}(\bar{s}_K, y)$  in DTRs.

**Theorem 4.** *Given  $P(\bar{s}_K, \bar{x}_K, y) > 0$ , there exists DTRs  $M_1, M_2$  such that  $P^{M_1}(\bar{s}_K, \bar{x}_K, y) = P^{M_2}(\bar{s}_K, \bar{x}_K, y) = P(\bar{s}_K, \bar{x}_K, y)$  while  $P_{\bar{x}_K}^{M_1}(\bar{s}_K, y) \neq P_{\bar{x}_K}^{M_2}(\bar{s}_K, y)$ .*

*Proof.* We define two counterfactual DTRs  $M_1, M_2 \in \mathcal{M}_{\text{OBS}}$  that are compatible with the observational distribution  $P(\bar{s}_K, \bar{x}_K, y)$ . If  $K = 1$ , for any  $y, s_1, x_1$  and any  $x_1^y \neq x_1$ , we define

$$P^{M_1}(y_{x_1^y} | s_1, x_1) = 0, \quad P^{M_2}(y_{x_1^y} | s_1, x_1) = 1$$

It is verifiable that

$$P_{x_1}^{M_1}(s_1, y) = P(s_1, x_1, y), \quad P_{x_1}^{M_2}(s_1, y) = P(s_1, x_1, y) + (1 - P(x_1 | s_1))P(s_1)$$

Since  $P(\bar{s}_K, \bar{x}_K, y) > 0$ , we have  $P_{x_1}^{M_2}(s_1, y) \neq P_{x_1}^{M_1}(s_1, y)$ .

We now consider the case where  $K > 1$ . For any  $\bar{x}_K, \bar{s}_K, y$ , and any  $\bar{x}_K^y \neq \bar{x}_K$ , we define

$$P^{M_1}(y_{\bar{x}_K^y} | \bar{s}_{K\bar{x}_{K-1}}, \bar{x}_K) = 0 \tag{33}$$

By definitions,  $P_{\bar{x}_K}^{M_1}(\bar{s}_K, y)$  is equal to the counterfactual quantities  $P^{M_1}(\bar{s}_{K\bar{x}_{K-1}}, y_{\bar{x}_K})$ . Thus,

$$\begin{aligned}
P_{\bar{x}_K}^{M_1}(\bar{s}_K, y) &= P^{M_1}(\bar{s}_{K\bar{x}_{K-1}}, y_{\bar{x}_K}, \bar{x}_K) + \sum_{\bar{x}_K' \neq \bar{x}_K} P^{M_1}(\bar{s}_{K\bar{x}_{K-1}}, y_{\bar{x}_K}, \bar{x}_K') \\
&= P^{M_1}(\bar{s}_{K\bar{x}_{K-1}}, y_{\bar{x}_K}, \bar{x}_K) + \sum_{\bar{x}_K' \neq \bar{x}_K} P^{M_1}(y_{\bar{x}_K} | \bar{s}_{K\bar{x}_{K-1}}, \bar{x}_K') P^{M_1}(\bar{s}_{K\bar{x}_{K-1}}, \bar{x}_K')
\end{aligned}$$

By the composition axiom,  $\bar{S}_{K\bar{x}_{K-1}} = \bar{S}_K, Y_{\bar{x}_K} = Y$  if  $\bar{X}_K = \bar{x}_K$ . Thus,  $P^{M_1}(\bar{s}_{K\bar{x}_{K-1}}, y_{\bar{x}_K}, \bar{x}_K) = P^{M_1}(\bar{s}_K, y, \bar{x}_K)$ . Since  $M_1 \in \mathcal{M}_{\text{OBS}}$ ,  $P^{M_1}(\bar{s}_K, y, \bar{x}_K) = P(\bar{s}_K, y, \bar{x}_K)$ . Together with Eq. (33), we can obtain

$$P_{\bar{x}_K}^{M_1}(\bar{s}_K, y) = P(\bar{s}_K, \bar{x}_K, y).$$

As for  $M_2$ , for any  $\bar{x}_{K-1}^K \neq \bar{x}_{K-1}$ , we define its factor

$$P^{M_2}(s_{K\bar{x}_{K-1}^K} | \bar{s}_{K-1\bar{x}_{K-2}}, \bar{x}_{K-1}) = 0$$

The above equation implies that for any  $\bar{x}'_{K-1} \neq \bar{x}_{K-1}$ ,

$$\begin{aligned} & P^{M_2}(\bar{s}_{K\bar{x}_{K-1}}, y_{\bar{x}_K}, \bar{x}'_{K-1}) \\ &= P^{M_2}(y_{\bar{x}_K} | \bar{s}_{K\bar{x}_{K-1}}, \bar{x}'_{K-1}) P^{M_2}(s_{K\bar{x}_{K-1}} | \bar{s}_{K-1\bar{x}_{K-2}}, \bar{x}'_{K-1}) P^{M_2}(\bar{s}_{K-1\bar{x}_{K-2}}, \bar{x}'_{K-1}) \\ &= 0 \end{aligned} \quad (34)$$

For any  $\bar{x}_K^y \neq \bar{x}_K$ , we define

$$P^{M_2}(y_{\bar{x}_K} | \bar{s}_{K\bar{x}_{K-1}}, \bar{x}_K) = 1 \quad (35)$$

We will now show that the above equation implies that for any  $x'_K \neq x_K$ ,

$$P^{M_2}(\bar{s}_{K\bar{x}_{K-1}}, y_{\bar{x}_K}, x'_K, \bar{x}_{K-1}) = P(\bar{s}_K, \bar{x}_{K-1}). \quad (36)$$

We first write  $P^{M_2}(\bar{s}_{K\bar{x}_{K-1}}, y_{\bar{x}_K}, x'_K, \bar{x}_{K-1})$  as:

$$P^{M_2}(\bar{s}_{K\bar{x}_{K-1}}, y_{\bar{x}_K}, x'_K, \bar{x}_{K-1}) = P^{M_2}(y_{\bar{x}_K} | \bar{s}_{K\bar{x}_{K-1}}, x'_K, \bar{x}_{K-1}) P^{M_2}(\bar{s}_{K\bar{x}_{K-1}}, x'_K, \bar{x}_{K-1})$$

It is immediate from Eq. (35) that

$$P^{M_2}(\bar{s}_{K\bar{x}_{K-1}}, y_{\bar{x}_K}, X_k \neq x_k, \bar{x}_{K-1}) = P^{M_2}(\bar{s}_{K\bar{x}_{K-1}}, \bar{x}_{K-1}).$$

By the composition axiom,  $\bar{S}_{K\bar{x}_{K-1}} = \bar{S}_K$  if  $\bar{X}_{K-1} = \bar{x}_{K-1}$ . Since  $M_2 \in \mathcal{M}_{\text{OBS}}$ , we thus have:

$$P^{M_2}(\bar{s}_{K\bar{x}_{K-1}}, y_{\bar{x}_K}, X_k \neq x_k, \bar{x}_{K-1}) = P^{M_2}(\bar{s}_K, \bar{x}_{K-1}) = P(\bar{s}_K, \bar{x}_{K-1}).$$

We now turn our attention to the interventional distribution  $P_{\bar{x}_K}^{M_2}(\bar{s}_K, y)$ . By expanding on  $\bar{X}_K$ ,

$$\begin{aligned} P_{\bar{x}_K}^{M_2}(\bar{s}_K, y) &= P^{M_2}(\bar{s}_{K\bar{x}_{K-1}}, y_{\bar{x}_K}, \bar{x}_K) + P^{M_2}(\bar{s}_{K\bar{x}_{K-1}}, y_{\bar{x}_K}, X_k \neq x_k, \bar{x}_{K-1}) \\ &\quad + P^{M_2}(\bar{s}_{K\bar{x}_{K-1}}, y_{\bar{x}_K}, \bar{X}_{K-1} \neq \bar{x}_{K-1}) \end{aligned}$$

The above equation, together with Eqs. (34) and (36), gives:

$$P_{\bar{x}_K}^{M_2}(\bar{s}_K, y) = P^{M_2}(\bar{s}_{K\bar{x}_{K-1}}, y_{\bar{x}_K}, \bar{x}_K) + P(\bar{s}_K, \bar{x}_{K-1}).$$

Again, by the composition axiom and  $M_2 \in \mathcal{M}_{\text{OBS}}$ ,

$$P_{\bar{x}_K}^{M_2}(\bar{s}_K, y) = P^{M_2}(\bar{s}_K, y, \bar{x}_K) + P(\bar{s}_K, \bar{x}_{K-1}) = P(\bar{s}_K, y, \bar{x}_K) + P(\bar{s}_K, \bar{x}_{K-1}).$$

Since  $P(\bar{s}_K, \bar{x}_{K-1}) > 0$ , we have  $P_{\bar{x}_K}^{M_1}(\bar{s}_K, y) \neq P_{\bar{x}_K}^{M_2}(\bar{s}_K, y)$ , which proves the statement.  $\square$

**Lemma 1.** For a DTR, given  $P(\bar{s}_K, \bar{x}_K, y)$ , for any  $k = 1, \dots, K-1$ ,

$$P_{\bar{x}_k}(\bar{s}_{k+1}) - P_{\bar{x}_k}(\bar{s}_k) \leq P(\bar{s}_{k+1}, \bar{x}_k) - P(\bar{s}_k, \bar{x}_k).$$

*Proof.* Note that  $P_{\bar{x}_k}(\bar{s}_{k+1})$  can be written as the counterfactual quantity  $P(\bar{s}_{k+1\bar{x}_k})$ . For any set of variables  $\mathbf{V}$ , let  $\neg v$  denote an event  $\mathbf{V} \neq v$ .  $P_{\bar{x}_k}(\bar{s}_{k+1})$  could thus be written as:

$$P_{\bar{x}_k}(\bar{s}_{k+1}) = P(\bar{s}_{k+1\bar{x}_k}, \bar{x}_k) + P(\bar{s}_{k+1\bar{x}_k}, \neg x_k, \bar{x}_{k-1}) + P(\bar{s}_{k+1\bar{x}_k}, \neg \bar{x}_{k-1}),$$

By the composition axiom,  $\bar{S}_{k+1\bar{x}_k} = \bar{S}_{k+1}$  if  $\bar{X}_k = \bar{x}_k$ . So,

$$\begin{aligned} P_{\bar{x}_k}(\bar{s}_{k+1}) &= P(\bar{s}_{k+1}, \bar{x}_k) + P(\bar{s}_{k+1\bar{x}_k}, \neg x_k, \bar{x}_{k-1}) + P(\bar{s}_{k+1\bar{x}_k}, \neg \bar{x}_{k-1}) \\ &\leq P(\bar{s}_{k+1}, \bar{x}_k) + P(\bar{s}_{k\bar{x}_k}, \neg x_k, \bar{x}_{k-1}) + P(\bar{s}_{k\bar{x}_k}, \neg \bar{x}_{k-1}) \\ &= P(\bar{s}_{k+1}, \bar{x}_k) + P(\bar{s}_{k\bar{x}_k}, \bar{x}_{k-1}) - P(\bar{s}_{k\bar{x}_k}, \bar{x}_k) + P(\bar{s}_{k\bar{x}_k}) - P(\bar{s}_{k\bar{x}_k}, \bar{x}_{k-1}) \\ &= P(\bar{s}_{k\bar{x}_k}) + P(\bar{s}_{k+1}, \bar{x}_k) - P(\bar{s}_{k\bar{x}_k}, \bar{x}_k). \end{aligned}$$

Again, by the composition axiom,  $\bar{S}_{k\bar{x}_k} = \bar{S}_k$  if  $\bar{X}_k = \bar{x}_k$ . Since  $P(\bar{s}_{k\bar{x}_k}) = P_{\bar{x}_k}(\bar{s}_k)$ ,

$$P_{\bar{x}_k}(\bar{s}_{k+1}) \leq P_{\bar{x}_k}(\bar{s}_k) + P(\bar{s}_{k+1}, \bar{x}_k) - P(\bar{s}_k, \bar{x}_k)$$

Rearranging the above equation proves the statement.  $\square$

**Lemma 4.** For a DTR, given  $P(\bar{s}_K, \bar{x}_K, y)$ , for any  $k = 0, \dots, K-1$ ,

$$P_{\bar{x}_k}(\bar{s}_{k+1}) \leq \Gamma(\bar{s}_{k+1}, \bar{x}_k),$$

where  $\Gamma(\bar{s}_{k+1}, \bar{x}_k) = P(\bar{s}_{k+1}, \bar{x}_k) - P(\bar{s}_k, \bar{x}_k) + \Gamma(\bar{s}_k, \bar{x}_{k-1})$  and  $\Gamma(s_1) = P(s_1)$ .

*Proof.* We prove this statement by induction.

**Base Case:**  $k = 0$  By definition,  $\Gamma(s_1) = P(s_1)$ . We thus have  $P(s_1) \leq \Gamma(s_1)$ .

**Induction Step** We assume that the statement holds for  $k$ , i.e.,  $P_{\bar{\mathbf{x}}_k}(\bar{\mathbf{s}}_{k+1}) \leq \Gamma(\bar{\mathbf{s}}_{k+1}, \bar{\mathbf{x}}_k)$ . We will prove that the statement holds for  $k + 1$ , i.e.,  $P_{\bar{\mathbf{x}}_{k+1}}(\bar{\mathbf{s}}_{k+2}) \leq \Gamma(\bar{\mathbf{s}}_{k+2}, \bar{\mathbf{x}}_{k+1})$ . To begin with,

$$P_{\bar{\mathbf{x}}_{k+1}}(\bar{\mathbf{s}}_{k+2}) = P_{\bar{\mathbf{x}}_{k+1}}(\bar{\mathbf{s}}_{k+2}) - P_{\bar{\mathbf{x}}_{k+1}}(\bar{\mathbf{s}}_{k+1}) + P_{\bar{\mathbf{x}}_{k+1}}(\bar{\mathbf{s}}_{k+1}).$$

By Lem. 1,

$$P_{\bar{\mathbf{x}}_{k+1}}(\bar{\mathbf{s}}_{k+2}) \leq P(\bar{\mathbf{s}}_{k+2}, \bar{\mathbf{x}}_{k+1}) - P(\bar{\mathbf{s}}_{k+1}, \bar{\mathbf{x}}_{k+1}) + P_{\bar{\mathbf{x}}_{k+1}}(\bar{\mathbf{s}}_{k+1}).$$

Since  $\bar{\mathbf{S}}_{k+1}$  are non-descendants of  $X_{k+1}$ ,  $P_{\bar{\mathbf{x}}_{k+1}}(\bar{\mathbf{s}}_{k+1}) = P_{\bar{\mathbf{x}}_k}(\bar{\mathbf{s}}_{k+1})$ . Since  $P_{\bar{\mathbf{x}}_k}(\bar{\mathbf{s}}_{k+1}) \leq \Gamma(\bar{\mathbf{s}}_{k+1}, \bar{\mathbf{x}}_k)$ ,

$$P_{\bar{\mathbf{x}}_{k+1}}(\bar{\mathbf{s}}_{k+2}) \leq P(\bar{\mathbf{s}}_{k+2}, \bar{\mathbf{x}}_{k+1}) - P(\bar{\mathbf{s}}_{k+1}, \bar{\mathbf{x}}_{k+1}) + \Gamma(\bar{\mathbf{s}}_{k+1}, \bar{\mathbf{x}}_k) = \Gamma(\bar{\mathbf{s}}_{k+2}, \bar{\mathbf{x}}_{k+1}). \quad \square$$

**Theorem 5.** For a DTR, given  $P(\bar{\mathbf{s}}_K, \bar{\mathbf{x}}_K, y) > 0$ , for any  $k = 1, \dots, K - 1$ ,

$$\frac{P(\bar{\mathbf{s}}_{k+1}, \bar{\mathbf{x}}_k)}{\Gamma(\bar{\mathbf{s}}_k, \bar{\mathbf{x}}_{k-1})} \leq P_{\bar{\mathbf{x}}_k}(s_{k+1}|\bar{\mathbf{s}}_k) \leq \frac{\Gamma(\bar{\mathbf{s}}_{k+1}, \bar{\mathbf{x}}_k)}{\Gamma(\bar{\mathbf{s}}_k, \bar{\mathbf{x}}_{k-1})},$$

*Proof.* By basic probabilistic operations,

$$P_{\bar{\mathbf{x}}_k}(s_{k+1}|\bar{\mathbf{s}}_k) = \frac{P_{\bar{\mathbf{x}}_k}(\bar{\mathbf{s}}_{k+1})}{P_{\bar{\mathbf{x}}_k}(\bar{\mathbf{s}}_k)}.$$

By Lem. 1,

$$P_{\bar{\mathbf{x}}_k}(s_{k+1}|\bar{\mathbf{s}}_k) \leq 1 + \frac{P(\bar{\mathbf{s}}_{k+1}, \bar{\mathbf{x}}_k) - P(\bar{\mathbf{s}}_k, \bar{\mathbf{x}}_k)}{P_{\bar{\mathbf{x}}_k}(\bar{\mathbf{s}}_k)}.$$

Since  $P(\bar{\mathbf{s}}_{k+1}, \bar{\mathbf{x}}_k) \leq P(\bar{\mathbf{s}}_k, \bar{\mathbf{x}}_k)$ ,  $P_{\bar{\mathbf{x}}_k}(s_{k+1}|\bar{\mathbf{s}}_k)$  is upper-bounded when  $P_{\bar{\mathbf{x}}_k}(\bar{\mathbf{s}}_k)$  is the maximal. Since  $\bar{\mathbf{S}}_k$  are non-descendants of  $X_k$ ,  $P_{\bar{\mathbf{x}}_k}(\bar{\mathbf{s}}_k) = P_{\bar{\mathbf{x}}_{k-1}}(\bar{\mathbf{s}}_k)$ . Together with Lem. 4, the above equation can be further bounded as:

$$P_{\bar{\mathbf{x}}_k}(s_{k+1}|\bar{\mathbf{s}}_k) \leq 1 + \frac{P(\bar{\mathbf{s}}_{k+1}, \bar{\mathbf{x}}_k) - P(\bar{\mathbf{s}}_k, \bar{\mathbf{x}}_k)}{\Gamma(\bar{\mathbf{s}}_k, \bar{\mathbf{x}}_{k-1})} = \frac{\Gamma(\bar{\mathbf{s}}_{k+1}, \bar{\mathbf{x}}_k)}{\Gamma(\bar{\mathbf{s}}_k, \bar{\mathbf{x}}_{k-1})}.$$

By definition,  $P_{\bar{\mathbf{x}}_k}(\bar{\mathbf{s}}_{k+1}) = P(\bar{\mathbf{s}}_{k+1}, \bar{\mathbf{x}}_k)$ . By basic probabilistic operations,

$$P_{\bar{\mathbf{x}}_k}(s_{k+1}|\bar{\mathbf{s}}_k) = \frac{P(\bar{\mathbf{s}}_{k+1}, \bar{\mathbf{x}}_k) + P(\bar{\mathbf{s}}_{k+1}, \neg \bar{\mathbf{x}}_k)}{P_{\bar{\mathbf{x}}_k}(\bar{\mathbf{s}}_k)} \geq \frac{P(\bar{\mathbf{s}}_{k+1}, \bar{\mathbf{x}}_k)}{P_{\bar{\mathbf{x}}_k}(\bar{\mathbf{s}}_k)}.$$

By the composition axiom,  $\bar{\mathbf{S}}_{k+1}, \bar{\mathbf{x}}_k = \bar{\mathbf{S}}_{k+1}$  if  $\bar{\mathbf{X}}_k = \bar{\mathbf{x}}_k$ . Applying Lem. 4 again gives

$$P_{\bar{\mathbf{x}}_k}(s_{k+1}|\bar{\mathbf{s}}_k) \geq \frac{P(\bar{\mathbf{s}}_{k+1}, \bar{\mathbf{x}}_k)}{P_{\bar{\mathbf{x}}_k}(\bar{\mathbf{s}}_k)} = \frac{P(\bar{\mathbf{s}}_{k+1}, \bar{\mathbf{x}}_k)}{\Gamma(\bar{\mathbf{s}}_k, \bar{\mathbf{x}}_{k-1})}. \quad \square$$

**Theorem 6.** Given  $P(\bar{\mathbf{s}}_K, \bar{\mathbf{x}}_K, y) > 0$ , for any  $k \in \{1, \dots, K - 1\}$ , let  $P_{\bar{\mathbf{x}}_k}(s_{k+1}|\bar{\mathbf{s}}_k) \in [a_{\bar{\mathbf{x}}_k, \bar{\mathbf{s}}_k}(s_{k+1}), b_{\bar{\mathbf{x}}_k, \bar{\mathbf{s}}_k}(s_{k+1})]$  denote the bound given by Thm. 5. There exists DTRs  $M_1, M_2$  such that  $P^{M_1}(\bar{\mathbf{s}}_K, \bar{\mathbf{x}}_K, y) = P^{M_2}(\bar{\mathbf{s}}_K, \bar{\mathbf{x}}_K, y) = P(\bar{\mathbf{s}}_K, \bar{\mathbf{x}}_K, y)$  while  $P_{\bar{\mathbf{x}}_k}^{M_1}(s_{k+1}|\bar{\mathbf{s}}_k) = a_{\bar{\mathbf{x}}_k, \bar{\mathbf{s}}_k}(s_{k+1})$ ,  $P_{\bar{\mathbf{x}}_k}^{M_2}(s_{k+1}|\bar{\mathbf{s}}_k) = b_{\bar{\mathbf{x}}_k, \bar{\mathbf{s}}_k}(s_{k+1})$ .

*Proof.* Without loss of generality, we assume that  $K > 1$ . We consider two counterfactual DTRs  $M_1, M_2 \in \mathcal{M}_{\text{OBS}}$  compatible with the observational distribution  $P(\bar{\mathbf{s}}_K, \bar{\mathbf{x}}_K, y)$ , which we define at the beginning of this section. For all  $i = 1, \dots, k - 1$ , for any  $\bar{\mathbf{x}}_i^{i+1} \neq \bar{\mathbf{x}}_i$ , we define that for any  $M \in \{M_1, M_2\}$ , its factors satisfy:

$$P^M(s_{i+1}|\bar{\mathbf{x}}_i^{i+1}, \bar{\mathbf{s}}_{i-1}, \bar{\mathbf{x}}_i) = 1. \quad (37)$$

Following a similar argument in Lem. 1, we will show that for any  $M \in \{M_1, M_2\}$ , for any  $i = 1, \dots, k - 1$ ,

$$P_{\bar{\mathbf{x}}_i}^M(\bar{\mathbf{s}}_{i+1}) - P_{\bar{\mathbf{x}}_i}^M(\bar{\mathbf{s}}_i) = P(\bar{\mathbf{s}}_{i+1}, \bar{\mathbf{x}}_i) - P(\bar{\mathbf{s}}_i, \bar{\mathbf{x}}_i). \quad (38)$$

By  $P_{\bar{\mathbf{x}}_i}^M(\bar{\mathbf{s}}_{i+1}) = P^M(\bar{\mathbf{s}}_{i+1|\bar{\mathbf{x}}_i})$  and basic probabilistic operations,

$$P_{\bar{\mathbf{x}}_i}^M(\bar{\mathbf{s}}_{i+1}) = P^M(\bar{\mathbf{s}}_{i+1|\bar{\mathbf{x}}_i}, \bar{\mathbf{x}}_i) + P^M(\bar{\mathbf{s}}_{i+1|\bar{\mathbf{x}}_i}, X_i \neq x_i, \bar{\mathbf{x}}_{i-1}) + P^M(\bar{\mathbf{s}}_{i+1|\bar{\mathbf{x}}_i}, \bar{X}_{i-1} \neq \bar{\mathbf{x}}_{i-1}).$$

By the composition axiom,  $\bar{\mathbf{S}}_{i+1|\bar{\mathbf{x}}_i} = \bar{\mathbf{S}}_{i+1}$  if  $\bar{X}_i = \bar{\mathbf{x}}_i$ . Since  $M \in \mathcal{M}_{\text{OBS}}$ ,  $P^M(\bar{\mathbf{s}}_{i+1}, \bar{\mathbf{x}}_i) = P(\bar{\mathbf{s}}_{i+1}, \bar{\mathbf{x}}_i)$ . Therefore,

$$\begin{aligned} P_{\bar{\mathbf{x}}_i}^M(\bar{\mathbf{s}}_{i+1}) &= P(\bar{\mathbf{s}}_{i+1}, \bar{\mathbf{x}}_i) + P^M(\bar{\mathbf{s}}_{i+1|\bar{\mathbf{x}}_i}, X_i \neq x_i, \bar{\mathbf{x}}_{i-1}) + P^M(\bar{\mathbf{s}}_{i+1|\bar{\mathbf{x}}_i}, \bar{X}_{i-1} \neq \bar{\mathbf{x}}_{i-1}), \\ &= P(\bar{\mathbf{s}}_{i+1}, \bar{\mathbf{x}}_i) + \sum_{x'_i \neq x_i} P^M(s_{i+1|\bar{\mathbf{x}}_i} | \bar{\mathbf{s}}_{i\bar{\mathbf{x}}_{i-1}}, x'_i, \bar{\mathbf{x}}_{i-1}) P(\bar{\mathbf{s}}_{i\bar{\mathbf{x}}_{i-1}}, x'_i, \bar{\mathbf{x}}_{i-1}) \\ &\quad + \sum_{\bar{\mathbf{x}}'_{i-1} \neq \bar{\mathbf{x}}_{i-1}} P^M(s_{i+1|\bar{\mathbf{x}}_i} | \bar{\mathbf{s}}_{i\bar{\mathbf{x}}_{i-1}}, x_i, \bar{\mathbf{x}}'_{i-1}) P(\bar{\mathbf{s}}_{i\bar{\mathbf{x}}_{i-1}}, x_i, \bar{\mathbf{x}}'_{i-1}) \end{aligned}$$

By Eq. (37),  $P^M(s_{i+1|\bar{\mathbf{x}}_i} | \bar{\mathbf{s}}_{i\bar{\mathbf{x}}_{i-1}}, x'_i, \bar{\mathbf{x}}_{i-1}) = P^M(s_{i+1|\bar{\mathbf{x}}_i} | \bar{\mathbf{s}}_{i\bar{\mathbf{x}}_{i-1}}, x_i, \bar{\mathbf{x}}'_{i-1}) = 1$ , which gives

$$\begin{aligned} P_{\bar{\mathbf{x}}_i}^M(\bar{\mathbf{s}}_{i+1}) &= P(\bar{\mathbf{s}}_{i+1}, \bar{\mathbf{x}}_i) + P^M(\bar{\mathbf{s}}_{i\bar{\mathbf{x}}_i}, X_i \neq x_i, \bar{\mathbf{x}}_{i-1}) + P^M(\bar{\mathbf{s}}_{i\bar{\mathbf{x}}_i}, \bar{X}_{i-1} \neq \bar{\mathbf{x}}_{i-1}) \\ &= P(\bar{\mathbf{s}}_{i+1}, \bar{\mathbf{x}}_i) + P^M(\bar{\mathbf{s}}_{i\bar{\mathbf{x}}_i}, \bar{\mathbf{x}}_{i-1}) - P^M(\bar{\mathbf{s}}_{i\bar{\mathbf{x}}_i}, \bar{\mathbf{x}}_i) + P^M(\bar{\mathbf{s}}_{i\bar{\mathbf{x}}_i}) - P^M(\bar{\mathbf{s}}_{i\bar{\mathbf{x}}_i}, \bar{\mathbf{x}}_{i-1}) \\ &= P^M(\bar{\mathbf{s}}_{i\bar{\mathbf{x}}_i}) + P(\bar{\mathbf{s}}_{i+1}, \bar{\mathbf{x}}_i) - P^M(\bar{\mathbf{s}}_{i\bar{\mathbf{x}}_i}, \bar{\mathbf{x}}_i) \end{aligned}$$

Again, by the composition axiom and  $M \in \mathcal{M}_{\text{OBS}}$ ,  $P^M(\bar{\mathbf{s}}_{i\bar{\mathbf{x}}_i}, \bar{\mathbf{x}}_i) = P(\bar{\mathbf{s}}_i, \bar{\mathbf{x}}_i)$ . Since  $P^M(\bar{\mathbf{s}}_{i\bar{\mathbf{x}}_i}) = P_{\bar{\mathbf{x}}_i}^M(\bar{\mathbf{s}}_i)$ , we have

$$P_{\bar{\mathbf{x}}_i}^M(\bar{\mathbf{s}}_{i+1}) = P_{\bar{\mathbf{x}}_i}^M(\bar{\mathbf{s}}_i) + P(\bar{\mathbf{s}}_{i+1}, \bar{\mathbf{x}}_i) - P(\bar{\mathbf{s}}_i, \bar{\mathbf{x}}_i).$$

Rearranging the above equation proves Eq. (37). Following a similar induction procedure in the proof of Lem. 4, we have that for any  $M \in \{M_1, M_2\}$ ,

$$P_{\bar{\mathbf{x}}_{k-1}}^M(\bar{\mathbf{s}}_k) = \Gamma(\bar{\mathbf{s}}_k, \bar{\mathbf{x}}_{k-1}). \quad (39)$$

As for  $M_1$ , for any  $\bar{\mathbf{x}}_k^{k+1} \neq \bar{\mathbf{x}}_k$ , we define

$$P^{M_1}(s_{k+1|\bar{\mathbf{x}}_k^{k+1}} | \bar{\mathbf{s}}_{k\bar{\mathbf{x}}_{k-1}}, \bar{\mathbf{x}}_k) = 0$$

This implies

$$\begin{aligned} P_{\bar{\mathbf{x}}_k}^{M_1}(\bar{\mathbf{s}}_{k+1}) &= P^{M_1}(\bar{\mathbf{s}}_{k+1|\bar{\mathbf{x}}_k}, \bar{\mathbf{x}}_k) + \sum_{\bar{\mathbf{x}}'_k \neq \bar{\mathbf{x}}_k} P^{M_1}(s_{k+1|\bar{\mathbf{x}}_k} | \bar{\mathbf{s}}_{k\bar{\mathbf{x}}_{k-1}}, \bar{\mathbf{x}}'_k) P^{M_1}(\bar{\mathbf{s}}_{k\bar{\mathbf{x}}_{k-1}}, \bar{\mathbf{x}}'_k) \\ &= P^{M_1}(\bar{\mathbf{s}}_{k+1|\bar{\mathbf{x}}_k}, \bar{\mathbf{x}}_k). \end{aligned}$$

By the composition axiom and  $M_1 \in \mathcal{M}_{\text{OBS}}$ ,  $P^{M_1}(\bar{\mathbf{s}}_{k+1|\bar{\mathbf{x}}_k}, \bar{\mathbf{x}}_k) = P(\bar{\mathbf{s}}_{k+1}, \bar{\mathbf{x}}_k)$ , which gives

$$P_{\bar{\mathbf{x}}_k}^{M_1}(\bar{\mathbf{s}}_{k+1}) = P(\bar{\mathbf{s}}_{k+1}, \bar{\mathbf{x}}_k).$$

The above equation, together with Eq. (39), gives:

$$P_{\bar{\mathbf{x}}_k}^{M_1}(s_{k+1} | \bar{\mathbf{s}}_k) = \frac{P_{\bar{\mathbf{x}}_k}^{M_1}(\bar{\mathbf{s}}_{k+1})}{P_{\bar{\mathbf{x}}_{k-1}}^M(\bar{\mathbf{s}}_k)} = \frac{P(\bar{\mathbf{s}}_{k+1}, \bar{\mathbf{x}}_k)}{\Gamma(\bar{\mathbf{s}}_k, \bar{\mathbf{x}}_{k-1})} = a_{\bar{\mathbf{x}}_k, \bar{\mathbf{s}}_k}(s_{k+1}).$$

As for  $M_2$ , for any  $\bar{\mathbf{x}}_k^{k+1} \neq \bar{\mathbf{x}}_k$ , we define

$$P^{M_2}(s_{k+1|\bar{\mathbf{x}}_k^{k+1}} | \bar{\mathbf{s}}_{k\bar{\mathbf{x}}_{k-1}}, \bar{\mathbf{x}}_k) = 1.$$

Following a similar procedure for proving Eq. (39), we have

$$P_{\bar{\mathbf{x}}_k}^M(\bar{\mathbf{s}}_{k+1}) = \Gamma(\bar{\mathbf{s}}_{k+1}, \bar{\mathbf{x}}_k).$$

Thus,

$$P_{\bar{\mathbf{x}}_k}^{M_2}(s_{k+1} | \bar{\mathbf{s}}_k) = \frac{P_{\bar{\mathbf{x}}_k}^{M_2}(\bar{\mathbf{s}}_{k+1})}{P_{\bar{\mathbf{x}}_{k-1}}^{M_2}(\bar{\mathbf{s}}_k)} = \frac{\Gamma(\bar{\mathbf{s}}_{k+1}, \bar{\mathbf{x}}_k)}{\Gamma(\bar{\mathbf{s}}_k, \bar{\mathbf{x}}_{k-1})} = b_{\bar{\mathbf{x}}_k, \bar{\mathbf{s}}_k}(s_{k+1}). \quad \square$$

**Corollary 2.** For a DTR, given  $P(\bar{s}_K, \bar{x}_K, y) > 0$ ,

$$\frac{E[Y|\bar{s}_K, \bar{x}_K]P(\bar{s}_K, \bar{x}_K)}{\Gamma(\bar{s}_K, \bar{x}_{K-1})} \leq E_{\bar{x}_K}[Y|\bar{s}_K] \leq 1 - \frac{(1 - E[Y|\bar{s}_K, \bar{x}_K])P(\bar{s}_K, \bar{x}_K)}{\Gamma(\bar{s}_K, \bar{x}_{K-1})}.$$

*Proof.* By basic probabilistic operations,

$$E_{\bar{x}_K}[Y|\bar{s}_K] = \frac{E_{\bar{x}_K}[Y|\bar{s}_K]P_{\bar{x}_K}(\bar{s}_K)}{P_{\bar{x}_K}(\bar{s}_K)}.$$

Note the counterfactual  $Y_{\bar{x}_K, \bar{s}_K}(\mathbf{u}) \in [0, 1]$ . Following a similar argument as Lem. 1,

$$E_{\bar{x}_K}[Y|\bar{s}_K]P_{\bar{x}_K}(\bar{s}_K) - P_{\bar{x}_K}(\bar{s}_K) \leq E[Y|\bar{s}_K, \bar{x}_K]P(\bar{s}_K, \bar{x}_K) - P(\bar{s}_K, \bar{x}_K).$$

This implies

$$E_{\bar{x}_K}[Y|\bar{s}_K] \leq 1 + \frac{(E[Y|\bar{s}_K, \bar{x}_K] - 1)P(\bar{s}_K, \bar{x}_K)}{P_{\bar{x}_K}(\bar{s}_K)}$$

Since  $E[Y|\bar{s}_K, \bar{x}_K] \leq 1$ ,  $E_{\bar{x}_K}[Y|\bar{s}_K]$  is upper-bounded when  $P_{\bar{x}_K}(\bar{s}_K)$  is the maximal. Since  $\bar{S}_K$  are non-descendants of  $X_K$ ,  $P_{\bar{x}_K}(\bar{s}_K) = P_{\bar{x}_{K-1}}(\bar{s}_K)$ . By Lem. 4,

$$E_{\bar{x}_K}[Y|\bar{s}_K] \leq 1 + \frac{(E[Y|\bar{s}_K, \bar{x}_K] - 1)P(\bar{s}_K, \bar{x}_K)}{\Gamma(\bar{s}_K, \bar{x}_{K-1})}.$$

By definition,  $P_{\bar{x}_K}(y, \bar{s}_K) = P(y_{\bar{x}_K}, \bar{s}_{K\bar{x}_{K-1}})$ . By basic probabilistic operations,

$$E_{\bar{x}_K}[Y|\bar{s}_K] \geq \frac{E[Y_{\bar{x}_K}|\bar{s}_{K\bar{x}_{K-1}}, \bar{x}_K]P(\bar{s}_{K\bar{x}_{K-1}}, \bar{x}_K)}{P_{\bar{x}_{K-1}}(\bar{s}_K)}.$$

By the composition axiom,  $\bar{S}_{K\bar{x}_{K-1}} = \bar{S}_{K-1}$ ,  $Y_{\bar{x}_K} = Y$  if  $\bar{X}_K = \bar{x}_K$ . Applying Lem. 4 gives

$$E_{\bar{x}_K}[Y|\bar{s}_K] \geq \frac{E[Y|\bar{s}_K, \bar{x}_K]P(\bar{s}_K, \bar{x}_K)}{P_{\bar{x}_K}(\bar{s}_K)} = \frac{E[Y|\bar{s}_K, \bar{x}_K]P(\bar{s}_K, \bar{x}_K)}{\Gamma(\bar{s}_K, \bar{x}_{K-1})}. \quad \square$$

## Proof of Theorems 7 and 8

**Lemma 5.** Fix  $\epsilon > 0$ ,  $\delta \in (0, 1)$ . With probability (w.p.) of at least  $1 - \delta$ , it holds for any  $T > 1$ ,  $R_\epsilon(T)$  of UC-DTR with parameter  $\delta$  and causal bounds  $\mathcal{C}$  is bounded by

$$R_\epsilon(T) \leq \min \left\{ 12K\sqrt{|\mathcal{S}||\mathcal{X}|T_\epsilon \log(2K|\mathcal{S}||\mathcal{X}|T/\delta)}, \|\mathcal{C}\|_1 T_\epsilon \right\} + 4K\sqrt{T_\epsilon \log(2T/\delta)}$$

*Proof.* Note that causal bounds  $\mathcal{C}$  is a set  $\{\mathcal{C}_1, \dots, \mathcal{C}_K\}$  where for  $k = 1, \dots, K-1$ ,

$$\begin{aligned} \mathcal{C}_k &= \left\{ \forall \bar{s}_{k+1}, \bar{x}_k : [a_{\bar{x}_k, \bar{s}_k}(s_{k+1}), b_{\bar{x}_k, \bar{s}_k}(s_{k+1})] \right\}, \\ \text{and } \mathcal{C}_K &= \left\{ \forall \bar{s}_K, \bar{x}_K : [a_{\bar{x}_K, \bar{s}_K}, b_{\bar{x}_K, \bar{s}_K}] \right\}. \end{aligned} \quad (40)$$

$\mathcal{M}^c$  is a set of DTRs such that for any  $M \in \mathcal{M}^c$ , its causal quantities  $P_{\bar{x}_k}(s_{k+1}|\bar{s}_k)$  and  $E_{\bar{x}_K}[Y|\bar{s}_K]$  satisfy the causal bounds  $\mathcal{C}$ , i.e.,

$$P_{\bar{x}_k}(s_{k+1}|\bar{s}_k) \in [a_{\bar{x}_k, \bar{s}_k}(s_{k+1}), b_{\bar{x}_k, \bar{s}_k}(s_{k+1})], \quad \text{and} \quad E_{\bar{x}_K}[Y|\bar{s}_K] \in [a_{\bar{x}_K, \bar{s}_K}, b_{\bar{x}_K, \bar{s}_K}]. \quad (41)$$

We assume that the causal bounds are always valid, i.e.,  $P(M^* \in \mathcal{M}^c) = 1$ . Let  $\mathcal{M}_t^c = \mathcal{M}_t \cap \mathcal{M}^c$ . By union bounds and Hoeffding's inequality (following a similar argument in [4, C.1]),

$$P(M^* \notin \mathcal{M}_t^c) \leq P(M^* \notin \mathcal{M}_t) \leq \frac{\delta}{4t^2}. \quad (42)$$

Since  $\sum_{t=1}^{\infty} \frac{1}{4t^2} \leq \frac{\pi^2}{24} \delta < \frac{\delta}{2}$ , it follows that with probability at least  $1 - \frac{\delta}{2}$ ,  $M^* \in \mathcal{M}_t^c$  for all episodes  $t = 1, 2, \dots$ .

Following the proof of Lem. 3, we have

$$R_\epsilon(T) \leq K\sqrt{6T_\epsilon \log(2T/\delta)} + \sqrt{\frac{3T_\epsilon \log(2T/\delta)}{2}} + \sum_{k=1}^{K-1} \sum_{t \in L_\epsilon} V_{\pi_t}(\bar{\mathbf{S}}_k^t, \bar{\mathbf{X}}_k^t; M_t^{(k)}) - V_{\pi_t}(\bar{\mathbf{S}}_k^t, \bar{\mathbf{X}}_k^t; M_t^{(k+1)}) \quad (43)$$

$$+ \sum_{t \in L_\epsilon} (E_{\bar{\mathbf{X}}_K^t}^{M_t} [Y | \bar{\mathbf{S}}_K^t] - E_{\bar{\mathbf{X}}_K^t} [Y | \bar{\mathbf{S}}_K^t]). \quad (44)$$

It thus suffices to bound quantities in Eqs. (43) and (44) separately.

**Bounding Eq. (43)** By Eq. (23) and basic probabilistic operations,

$$\begin{aligned} & V_{\pi_t}(\bar{\mathbf{S}}_k^t, \bar{\mathbf{X}}_k^t; M_t^{(k)}) - V_{\pi_t}(\bar{\mathbf{S}}_k^t, \bar{\mathbf{X}}_k^t; M_t^{(k+1)}) \\ &= \sum_{s_{k+1}} (P^{M_t}(s_{k+1} | \bar{\mathbf{S}}_k, \bar{\mathbf{X}}_k) - P(s_{k+1} | \bar{\mathbf{S}}_k, \bar{\mathbf{X}}_k)) V_{\pi_t}(s_{k+1}, \bar{\mathbf{S}}_k^t, \bar{\mathbf{X}}_k^t; M_t) \\ &\leq \left\| P_{\bar{\mathbf{x}}_k}^{M_t}(\cdot | \bar{\mathbf{s}}_k) - P_{\bar{\mathbf{x}}_k}(\cdot | \bar{\mathbf{s}}_k) \right\|_1 \max_{s_{k+1}} V_{\pi_t}(s_{k+1}, \bar{\mathbf{S}}_k^t, \bar{\mathbf{X}}_k^t; M_t) \\ &\leq \min \left\{ 2\sqrt{6|\mathcal{S}_{k+1}| \log(2K|\bar{\mathbf{S}}_k||\bar{\mathbf{X}}_k|T/\delta)} \frac{1}{\sqrt{\max\{1, N^t(\bar{\mathbf{S}}_k^t, \bar{\mathbf{X}}_k^t)\}}}, \|\mathbf{c}_k\|_1 \right\} \end{aligned}$$

The last step follows from Eqs. (17) and (41). We thus have

$$\begin{aligned} & \sum_{t \in L_\epsilon} V_{\pi_t}(\bar{\mathbf{S}}_k^t, \bar{\mathbf{X}}_k^t; M_t^{(k)}) - V_{\pi_t}(\bar{\mathbf{S}}_k^t, \bar{\mathbf{X}}_k^t; M_t^{(k+1)}) \\ &\leq \sum_{t \in L_\epsilon} \min \left\{ 2\sqrt{6|\mathcal{S}_{k+1}| \log(2K|\bar{\mathbf{S}}_k||\bar{\mathbf{X}}_k|T/\delta)} \frac{1}{\sqrt{\max\{1, N^t(\bar{\mathbf{S}}_k^t, \bar{\mathbf{X}}_k^t)\}}}, \|\mathbf{c}_k\|_1 \right\} \\ &\leq \min \left\{ \sum_{t \in L_\epsilon} 2\sqrt{6|\mathcal{S}_{k+1}| \log(2K|\bar{\mathbf{S}}_k||\bar{\mathbf{X}}_k|T/\delta)} \frac{1}{\sqrt{\max\{1, N^t(\bar{\mathbf{S}}_k^t, \bar{\mathbf{X}}_k^t)\}}}, \sum_{t \in L_\epsilon} \|\mathbf{c}_k\|_1 \right\} \\ &\leq \min \left\{ 2(\sqrt{2} + 1)\sqrt{6T_\epsilon|\bar{\mathbf{S}}_{k+1}||\bar{\mathbf{X}}_k| \log(2K|\bar{\mathbf{S}}_k||\bar{\mathbf{X}}_k|T/\delta)}, \|\mathbf{c}_k\|_1 T_\epsilon \right\} \end{aligned}$$

The last step follows from results in [4, D] and  $|L_\epsilon| = T_\epsilon$ . Eq. (43) could thus be written as:

$$\begin{aligned} & \sum_{k=1}^{K-1} \sum_{t \in L_\epsilon} V_{\pi_t}(\bar{\mathbf{S}}_k^t, \bar{\mathbf{X}}_k^t; M_t^{(k)}) - V_{\pi_t}(\bar{\mathbf{S}}_k^t, \bar{\mathbf{X}}_k^t; M_t^{(k+1)}) \\ &\leq \sum_{k=1}^{K-1} \min \left\{ 2(\sqrt{2} + 1)\sqrt{6T_\epsilon|\bar{\mathbf{S}}_{k+1}||\bar{\mathbf{X}}_k| \log(2K|\bar{\mathbf{S}}_k||\bar{\mathbf{X}}_k|T/\delta)}, \|\mathbf{c}_k\|_1 T_\epsilon \right\} \\ &\leq \min \left\{ \sum_{k=1}^{K-1} 2(\sqrt{2} + 1)\sqrt{6T_\epsilon|\bar{\mathbf{S}}_{k+1}||\bar{\mathbf{X}}_k| \log(2K|\bar{\mathbf{S}}_k||\bar{\mathbf{X}}_k|T/\delta)}, \sum_{k=1}^{K-1} \|\mathbf{c}_k\|_1 T_\epsilon \right\} \end{aligned}$$

Thus,

$$\begin{aligned} & \sum_{k=1}^{K-1} \sum_{t \in L_\epsilon} V_{\pi_t}(\bar{\mathbf{S}}_k^t, \bar{\mathbf{X}}_k^t; M_t^{(k)}) - V_{\pi_t}(\bar{\mathbf{S}}_k^t, \bar{\mathbf{X}}_k^t; M_t^{(k+1)}) \\ &\leq \min \left\{ (K-1)2(\sqrt{2} + 1)\sqrt{6T_\epsilon|\bar{\mathbf{S}}||\bar{\mathbf{X}}| \log(2K|\bar{\mathbf{S}}||\bar{\mathbf{X}}|T/\delta)}, \sum_{k=1}^{K-1} \|\mathbf{c}_k\|_1 T_\epsilon \right\}. \end{aligned} \quad (45)$$

**Bounding Eq. (44)** Since both  $M^*, M_t$  are in the set  $\mathcal{M}_t^c$ ,

$$\begin{aligned} & E_{\bar{\mathbf{X}}_K^t}^{M_t} [Y | \bar{\mathbf{S}}_K^t] - E_{\bar{\mathbf{X}}_K^t} [Y | \bar{\mathbf{S}}_K^t] \leq \left| E_{\bar{\mathbf{x}}_K}^{M_t} [Y | \bar{\mathbf{s}}_K] - \hat{E}_{\bar{\mathbf{x}}_K}^t [Y | \bar{\mathbf{s}}_K] \right| + \left| E_{\bar{\mathbf{X}}_K^t} [Y | \bar{\mathbf{S}}_K^t] - \hat{E}_{\bar{\mathbf{x}}_K}^t [Y | \bar{\mathbf{s}}_K] \right| \\ &\leq \min \left\{ 2\sqrt{2 \log(2K|\bar{\mathbf{S}}||\bar{\mathbf{X}}|T/\delta)} \frac{1}{\sqrt{\max\{1, N^t(\bar{\mathbf{S}}_K^t, \bar{\mathbf{X}}_K^t)\}}}, \|\mathbf{c}_K\|_1 \right\} \end{aligned}$$

Eq. (44) can thus be written as:

$$\begin{aligned}
& \sum_{t \in L_\epsilon} (E_{\bar{\mathbf{X}}_K^t}^{M_t} [Y | \bar{\mathbf{S}}_K^t] - E_{\bar{\mathbf{X}}_K^t} [Y | \bar{\mathbf{S}}_K^t]) \\
& \leq \sum_{t \in L_\epsilon} \min \left\{ 2\sqrt{2 \log(2K|\mathcal{S}||\mathcal{X}|T/\delta)} \frac{1}{\sqrt{\max\{1, N^t(\bar{\mathbf{S}}_K^t, \bar{\mathbf{X}}_K^t)\}}}, \|\mathbf{c}_K\|_1 \right\} \\
& \leq \min \left\{ \sum_{t \in L_\epsilon} 2\sqrt{2 \log(2K|\mathcal{S}||\mathcal{X}|T/\delta)} \frac{1}{\sqrt{\max\{1, N^t(\bar{\mathbf{S}}_K^t, \bar{\mathbf{X}}_K^t)\}}}, \sum_{t \in L_\epsilon} \|\mathbf{c}_K\|_1 \right\}.
\end{aligned}$$

The last step follows from Eqs. (18) and (41). From results in [4, D], we have

$$\begin{aligned}
& \sum_{t \in L_\epsilon} (E_{\bar{\mathbf{X}}_K^t}^{M_t} [Y | \bar{\mathbf{S}}_K^t] - E_{\bar{\mathbf{X}}_K^t} [Y | \bar{\mathbf{S}}_K^t]) \\
& \leq \min \left\{ 2(\sqrt{2} + 1) \sqrt{2T_\epsilon |\bar{\mathcal{S}}| |\bar{\mathcal{X}}| \log(2K|\bar{\mathcal{S}}| |\bar{\mathcal{X}}| T/\delta)}, \|\mathbf{c}_K\|_1 T_\epsilon \right\}.
\end{aligned} \tag{46}$$

Eqs. (45) and (46) together give:

$$\begin{aligned}
R_\epsilon(T) & \leq K \sqrt{6T_\epsilon \log(2T/\delta)} + \sqrt{\frac{3T_\epsilon \log(2T/\delta)}{2}} \\
& + \min \left\{ (K-1)2(\sqrt{2} + 1) \sqrt{6T_\epsilon |\bar{\mathcal{S}}| |\bar{\mathcal{X}}| \log(2K|\bar{\mathcal{S}}| |\bar{\mathcal{X}}| T/\delta)}, \sum_{k=1}^{K-1} \|\mathbf{c}_k\|_1 T_\epsilon \right\} \\
& + \min \left\{ 2(\sqrt{2} + 1) \sqrt{2T_\epsilon |\bar{\mathcal{S}}| |\bar{\mathcal{X}}| \log(2K|\bar{\mathcal{S}}| |\bar{\mathcal{X}}| T/\delta)}, \|\mathbf{c}_K\|_1 T_\epsilon \right\}.
\end{aligned} \tag{47}$$

A quick simplification gives:

$$R_\epsilon(T) \leq \min \left\{ 12K \sqrt{|\mathcal{S}||\mathcal{X}| T_\epsilon \log(2K|\mathcal{S}||\mathcal{X}| T/\delta)}, \|\mathbf{c}\|_1 T_\epsilon \right\} + 4K \sqrt{T_\epsilon \log(2T/\delta)}. \quad \square$$

**Theorem 7.** Fix a  $\delta \in (0, 1)$ . With probability of at least  $1 - \delta$ , it holds for any  $T > 1$ , the regret of  $UC^C$ -DTR with parameter  $\delta$  and causal bounds  $\mathcal{C}$  is bounded by

$$R(T) \leq \min \left\{ 12K \sqrt{|\mathcal{S}||\mathcal{X}| T \log(2K|\mathcal{S}||\mathcal{X}| T/\delta)}, \|\mathbf{c}\|_1 T \right\} + 4K \sqrt{T \log(2T/\delta)}.$$

*Proof.* Fix  $\epsilon = 0$ . Naturally,  $T_\epsilon = T$  and  $R_\epsilon(T) = R(T)$ . By Lem. 5,

$$R(T) \leq \min \left\{ 12K \sqrt{|\mathcal{S}||\mathcal{X}| T \log(2K|\mathcal{S}||\mathcal{X}| T/\delta)}, \|\mathbf{c}\|_1 T \right\} + 4K \sqrt{T \log(2T/\delta)}. \quad \square$$

**Theorem 8.** For any  $T \geq 1$ , with parameter  $\delta = \frac{1}{T}$  and causal bounds  $\mathcal{C}$ , the expected regret of  $UC^C$ -DTR is bounded by

$$E[R(T)] \leq \max_{\pi \in \Pi_{\mathcal{C}}} \left\{ \frac{33^2 K^2 |\mathcal{S}||\mathcal{X}| \log(T)}{\Delta_\pi} + \frac{32}{\Delta_\pi^3} + \frac{4}{\Delta_\pi} \right\} + 1.$$

*Proof.* Let  $\tilde{R}_\epsilon(T)$  denote the regret cumulated in  $\epsilon$ -good episode up to  $T$  steps. By Eqs. (42) and (47),

$$\begin{aligned}
E[R(T)] & \leq E[R_\epsilon(T) I_{M^* \in \mathcal{M}_t^c}] + E[\tilde{R}_\epsilon(T) I_{M^* \in \mathcal{M}_t^c}] + \sum_{t=1}^T P(M \notin \mathcal{M}_t^c) \\
& \leq \min \left\{ 12K \sqrt{|\mathcal{S}||\mathcal{X}| T \log(2K|\mathcal{S}||\mathcal{X}| T/\delta)}, \|\mathbf{c}\|_1 T \right\} + 4K \sqrt{T \log(2T/\delta)} \\
& + E[\tilde{R}_\epsilon(T) I_{M^* \in \mathcal{M}_t^c}] + \frac{\delta}{T} \\
& \leq 23K \sqrt{|\mathcal{S}||\mathcal{X}| T_\epsilon \log(T/\delta)} + E[\tilde{R}_\epsilon(T) I_{M^* \in \mathcal{M}_t^c}] + \frac{\delta}{T}
\end{aligned}$$



Fix  $\delta = \frac{1}{T}$ , it is immediate from Eq. (29) that

$$E[R(T)] \leq \frac{23^2 K^2 |\mathcal{S}| |\mathcal{X}| \log(T^2)}{\epsilon} + E[\tilde{R}_\epsilon(T) I_{M^* \in \mathcal{M}_t^c}] + 1. \quad (48)$$

Note that when  $M^* \in \mathcal{M}_t^c$ , the maximal expected reward of any  $\pi_t$  over all instances in the family of DTRs  $\mathcal{M}_t^c$  must be no less than the true optimal value  $V_{\pi^*}(M^*)$ . In words,  $\Pi_G^-$  is the effective policy space of UC<sup>c</sup>-DTR procedure. Let  $\Delta = \arg \min_{\pi \in \Pi_G^-} \Delta_\pi$ . Fix  $\epsilon = \frac{\Delta}{2}$ , Eq. (48) implies:

$$E[R(T)] \leq \frac{33^2 K^2 |\mathcal{S}| |\mathcal{X}| \log(T)}{\Delta} + E[\tilde{R}_{\frac{\Delta}{2}}(T) I_{M^* \in \mathcal{M}_t^c}] + 1.$$

Among quantities in the above equation,  $E[\tilde{R}_{\frac{\Delta}{2}}(T) I_{M^* \in \mathcal{M}_t^c}]$  can be bounded following a similar procedure in the proof of Thm. 2, which proves the statement.  $\square$

## Appendix II. Experimental Setup

In this section, we provide details about the setup of experiments in the main text. For all experiments, we test sequentially randomized trials (*rand*), UC-DTR algorithm (*uc-dtr*) and the causal UC-DTR (*uc<sup>c</sup>-dtr*) with causal bounds derived from  $1 \times 10^5$  observational samples. Each experiment lasts for  $T = 1.1 \times 10^4$  episodes. The parameter  $\delta = 1/KT$  for *uc-dtr* and *uc<sup>c</sup>-dtr* where  $K$  is the total stages of interventions. For all algorithms, we measure their cumulative regret over 200 repetitions.

**Cancer Treatment** We test the survival model of the two-stage clinical trial conducted by the Cancer and Leukemia Group B [5, 8]. Protocol 8923 was a double-blind, placebo controlled two-stage trial reported by [7] examining the effects of infusions of granulocyte-macrophage colony-stimulating factor (GM-CSF) after initial chemotherapy in patients with acute myelogenous leukemia (AML). Standard chemotherapy for AML could place patients at increased risk of death due to infection or bleeding-related complications. GM-CSF administered after chemotherapy might assist patient recovery, thus reducing the number of deaths due to such complications. Patients were randomized initially to GM-CSF or placebo following standard chemotherapy. Later, patients meeting the criteria of complete remission and consenting to further participation were offered a second randomization to one of two intensification treatments.

We will describe this treatment procedure using the DTR with  $K = 2$ .  $X_1, X_2 \in \{0, 1\}$  represent treatments;  $S_1 = \emptyset$  and  $S_2$  indicates the observed remission after the first treatment (0 stands for no remission and 1 for complete remission);  $Y$  indicates the survival of patients at the time of recording. The exogenous variable  $U$  is the age of patients where  $U = 1$  if the patient is old and  $U = 0$  otherwise. Values of  $U$  are drawn from a distribution  $P(u)$  where  $P(U = 1) = 0.2358$ . Values of  $S_2$  are drawn from a distribution  $P_{x_1}(s_2)$  described in Table 1.

|         | $X_1 = 0$ | $X_1 = 1$ |
|---------|-----------|-----------|
| $U = 0$ | 0.8101    | 0.0883    |
| $U = 1$ | 0.7665    | 0.2899    |

Table 1: Probabilities of the distribution  $P(S_2 = 1|u, x_1)$ .

Let  $T_1, T_2$  denote the potential survival time induced by treatment  $X_1, X_2$  respectively. Values of  $T_1, T_2$  are decided by functions defined as follows:

$$T_1 \leftarrow \min\{(1 - S_2)T_1^* + S_2(T_2^* + T_3^*), L\}, \quad T_2 \leftarrow \min\{(1 - S_2)T_1^* + S_2(T_2^* + T_4^*), L\}$$

where  $L = 1.5$ . Let  $\exp(\beta)$  denote an exponential distribution with mean  $1/\beta$ . Values of  $T_1^*, T_2^*, T_3^*$  are drawn from exponential distributions defined as follows:

$$T_1^* \sim \exp(\beta_{u,x_1}^1), \quad T_2^* \sim \exp(\beta_{u,x_1}^2), \quad T_3^* \sim \exp(\beta_{u,x_1}^3)$$

Given  $T_3^*$ , values of  $T_4^*$  are drawn from distribution

$$T_4^* \sim \exp(\beta_{u,x_1}^3 + \beta_{u,x_1}^4 T_3^*).$$

The total survival time  $T$  of a patient is decided as follows:

$$T \leftarrow (1 - S_2)T_1 + S_2(1 - X_2)T_1 + S_2X_2T_2.$$

The parameters  $\beta_{u,x_1} = (\beta_{u,x_1}^1, \beta_{u,x_1}^2, \beta_{u,x_1}^3, \beta_{u,x_1}^4)$  are described in Table 2.

|         |           | $\beta_{u,x_1}^1$ | $\beta_{u,x_1}^2$ | $\beta_{u,x_1}^3$ | $\beta_{u,x_1}^4$ |
|---------|-----------|-------------------|-------------------|-------------------|-------------------|
| $U = 0$ | $X_1 = 0$ | 4.3063            | 4.9607            | 0.8737            | 4.2538            |
|         | $X_1 = 1$ | 0.8286            | 8.2074            | 8.7975            | 7.6468            |
| $U = 1$ | $X_1 = 0$ | 2.6989            | 0.0235            | 5.9835            | 6.8059            |
|         | $X_1 = 1$ | 3.6036            | 1.1007            | 9.4426            | 7.3960            |

Table 2: Parameters  $\beta_{u,x_1}$ .

The primary outcome  $Y$  is the survival of the patient at the time of observation  $t = 1$ . Values of  $Y$  are decided by the indicator function  $Y \leftarrow I_{T>1}$ .

We generate the confounded observational data following a sequence of decision rules  $X_1 \sim \pi_1(X_1|U)$ ,  $X_2 \sim \pi_2(X_2|U, X_1, S_2)$ . The policy  $\pi_1(X_1|U)$  is a conditional distribution mapping from  $U$  to the domain of  $X_1$  where  $\pi_1(X_1 = 1|U = 0) = 0.5102$  and  $\pi_1(X_1 = 1|U = 1) = 0.2433$ . Similarly,  $\pi_2(X_2|U, X_1, S_2)$  is a conditional distribution mapping from  $U, X_1, S_2$  to the domain of  $X_2$ ; Table 3 describes its parametrization.

|         | $X_1 = 0$ |           | $X_1 = 1$ |           |
|---------|-----------|-----------|-----------|-----------|
|         | $S_2 = 0$ | $S_2 = 1$ | $S_2 = 0$ | $S_2 = 1$ |
| $U = 0$ | 0.2173    | 0.8696    | 0.6195    | 0.4641    |
| $U = 1$ | 0.8869    | 0.0103    | 0.5314    | 0.4339    |

Table 3: Probabilities of  $\pi_2(X_2 = 1|U, X_1, S_2)$ .

## References

- [1] P. Auer, N. Cesa-Bianchi, Y. Freund, and R. E. Schapire. The nonstochastic multiarmed bandit problem. *SIAM journal on computing*, 32(1):48–77, 2002.
- [2] A. Charnes and W. W. Cooper. Programming with linear fractional functionals. *Naval Research logistics quarterly*, 9(3-4):181–186, 1962.
- [3] W. HOEFFDING. Probability inequalities for sums of bounded random variables. *J. Am. Stat. Associ.*, 58(301):13–30, 1963.
- [4] T. Jaksch, R. Ortner, and P. Auer. Near-optimal regret bounds for reinforcement learning. *Journal of Machine Learning Research*, 11(Apr):1563–1600, 2010.
- [5] J. K. Lunceford, M. Davidian, and A. A. Tsiatis. Estimation of survival distributions of treatment policies in two-stage randomization designs in clinical trials. *Biometrics*, 58(1):48–57, 2002.
- [6] J. Pearl. *Causality: Models, Reasoning, and Inference*. Cambridge University Press, New York, 2000. 2nd edition, 2009.
- [7] R. M. Stone, D. T. Berg, S. L. George, R. K. Dodge, P. A. Paciucci, P. Schulman, E. J. Lee, J. O. Moore, B. L. Powell, and C. A. Schiffer. Granulocyte–macrophage colony-stimulating factor after initial chemotherapy for elderly patients with primary acute myelogenous leukemia. *New England Journal of Medicine*, 332(25):1671–1677, 1995.
- [8] A. S. Wahed and A. A. Tsiatis. Optimal estimator for the survival distribution and related quantities for treatment policies in two-stage randomization designs in clinical trials. *Biometrics*, 60(1):124–133, 2004.
- [9] T. Weissman, E. Ordentlich, G. Seroussi, S. Verdu, and M. J. Weinberger. Inequalities for the 11 deviation of the empirical distribution. 2003.