

# ELIGIBILITY TRACES FOR CONFOUNDING ROBUST OFF-POLICY EVALUATION: A CAUSAL APPROACH

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

A unifying theme in Artificial Intelligence is learning an effective policy to control an agent in an unknown environment in order to optimize a certain performance measure. Off-policy methods can significantly improve the sample efficiency during training since they allow an agent to learn from observed trajectories generated by different *behavior policies*, without directly deploying the *target policies* in the underlying environment. This paper studies off-policy evaluation from biased offline data where (1) *unobserved confounding* bias cannot be ruled out a priori; or (2) the observed trajectories do not *overlap* with intended behaviors of the learner, i.e., the target and behavior policies do not share a common support. Specifically, we first extend the Bellman’s equation to derive effective closed-form bounds over value functions from the observational distribution contaminated with unobserved confounding and no-overlap. Second, we propose two novel algorithms that use eligibility traces to estimate these bounds from finite observational data. Compared to other partial identification methods for off-policy evaluation in sequential environments, these methods are model-free and do not rely on additional parametric knowledge about the system dynamics in the underlying environment.

## 1 INTRODUCTION

A typical reinforcement learning agent learns from past data, i.e., from observed trajectories of states, actions, and reward signals generated by the agent intervening in the underlying environment. This data reflects the influence of the decision-making policy used to allocate actions based on the observed state, which is called the *behavior policy*. This policy might be selected by the agent in the past or by a different demonstrator operating in the same environment. *Policy evaluation* studies the problem of evaluating the effectiveness of a candidate *target policy* from the combination of past data and theoretical assumptions about the environment. When the behavior and target policies coincide, the evaluation is called *on-policy* learning, in which the expected return of candidate policies given the agent’s starting state (i.e., the value function) could be directly estimated with empirical means (Sutton & Barto, 1998). In practice, however, the learner might have to learn about policies different from the currently deployed one that generated the data, leading to the *off-policy* learning problem.

Off-policy learning is a popular area of research, as it allows for more efficient learning by using data from different policies. Several algorithms have been proposed for off-policy evaluation from finite observations, including Q-learning (Watkins, 1989; Watkins & Dayan, 1992), importance sampling (Swaminathan & Joachims, 2015; Jiang & Li, 2016), and temporal difference (Precup et al., 2000; Munos et al., 2016). These algorithms rely on two critical assumptions about the behavior policy. First, no unobserved confounder affects the behavior policy’s selected action and the subsequent state and reward. Second, the behavior policy is stochastic, covering all intended actions the target policy selects given all observed states. When either of these assumptions does not hold, the effect of the target policy is generally not *identifiable*, i.e., the model assumptions are insufficient to uniquely determine the value function from the offline data (Pearl, 2000; Zhang & Bareinboim, 2019).

In recent times, researchers have been using partial identification methods to obtain reliable off-policy evaluation in situations where there are unobserved confounders, and the behavior and target policies have no common support (Kallus & Zhou, 2018; Zhang & Bareinboim, 2019; Kallus & Zhou, 2020; Namkoong et al., 2020; Khan et al., 2023; Bruns-Smith & Zhou, 2023; Kausik et al., 2024). Partial identification is a well-studied problem in causal inference (Balke & Pearl, 1997; Zhang et al., 2022),

054  
055  
056  
057  
058  
059  
060  
061  
062  
063  
064  
065  
066  
067  
068  
069  
070  
071  
072  
073  
074  
075  
076  
077  
078  
079  
080  
081  
082  
083  
084  
085  
086  
087  
088  
089  
090  
091  
092  
093  
094  
095  
096  
097  
098  
099  
100  
101  
102  
103  
104  
105  
106  
107

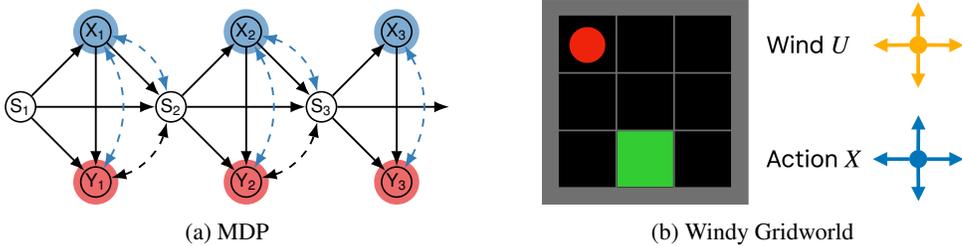


Figure 1: (a) Causal diagram representing the data-generating mechanisms in a Markov Decision Process (MDP); (b) A windy gridworld environment where the red dot represents the agent and green square is the goal state; the agent can take five actions - up, down, right, left, and stay-put; the wind can blow in five direction - north, south, west, east, and no-wind.

econometrics (Imbens & Rubin, 1997; Poirier, 1998; Romano & Shaikh, 2008; Stoye, 2009; Bugni, 2010; Todem et al., 2010; Moon & Schorfheide, 2012), and dynamical systems (Bajari et al., 2007; Norets & Tang, 2014; Dickstein & Morales, 2018; Morales et al., 2019; Berry & Compiani, 2023). It enables the derivation of informative bounds on target effects from confounded observational data. Among these works, researchers often employ a combination of approaches and constraints. These include (1) the marginal sensitivity model presuming access to a bound over the odds ratio between the nominal and actual behavioral policies (Kallus & Zhou, 2018; 2020; Namkoong et al., 2020; Khan et al., 2023; Bruns-Smith & Zhou, 2023); (2) additional parametric assumptions about the system dynamics (i.e., reward function and transition distribution) are invoked under which bounds are derived (Khan et al., 2023; Kausik et al., 2024); (3) a model-based algorithm is applied, which requires estimation of the underlying system dynamics (Zhang & Bareinboim, 2019); (4) the decision horizon is finite, i.e., the agent only determines a finite number of actions (Kallus & Zhou, 2018; Zhang & Bareinboim, 2019; Namkoong et al., 2020; Khan et al., 2023; Kausik et al., 2024). See Appendix A for a more detailed survey on partial identification and robust reinforcement learning.

This paper studies model-free algorithms for robust off-policy evaluation from confounded offline data generated by behavior policy with no-overlap support. We propose novel partial identification algorithms using eligibility traces to obtain informative bounds over the expected return of candidate policies from offline data generated from an unknown Markov decision process with an infinite horizon. More specifically, our contributions are summarized as follows. (1) We extend the Bellman equation that permits one to derive optimal bounds over target value functions from the observational distribution generated by an unknown behavior policy. (2) We propose a causal off-policy temporal difference algorithm (C-TD( $\lambda$ )) using eligibility traces to estimate bounds over the state value function from finite observations contaminated with unobserved confounding and no-overlap. (3) We introduce an alternative eligibility traces algorithm following tree backup (C-TB( $\lambda$ )) that obtains bounds over the state-action value function from confounded observations. Finally, we evaluate our proposed algorithms using extensive simulations in synthetic environments. Due to space constraints, all proofs are provided in Appendix B; details on the experiment setup are provided in Appendix C.

**Preliminaries and Notations** We use capital letters to denote random variables ( $X$ ), small letters for their values ( $x$ ) and  $\mathcal{D}_X$  for the domain of  $X$ . For an arbitrary set  $\mathbf{X}$ , let  $|\mathbf{X}|$  be its cardinality. Fix indices  $i, j \in \mathbb{N}$ . Let  $\mathbf{X}_{i:j}$  stand for a sequence of variables  $\{X_i, X_{i+1}, \dots, X_j\}$ ;  $\mathbf{X}_{i:j} = \emptyset$  if  $j < i$ . We denote by  $P(\mathbf{X})$  represents a probability distribution over variables  $\mathbf{X}$ . Similarly,  $P(\mathbf{Y} | \mathbf{X})$  represents a set of conditional distributions  $P(\mathbf{Y} | \mathbf{X} = \mathbf{x})$  for all realizations  $\mathbf{x}$ . We will consistently use  $P(\mathbf{x})$  as abbreviations for probabilities  $P(\mathbf{X} = \mathbf{x})$ ; so does  $P(\mathbf{Y} = \mathbf{y} | \mathbf{X} = \mathbf{x}) = P(\mathbf{y} | \mathbf{x})$ . Finally,  $\mathbb{1}_{Z=z}$  is an indicator function that returns 1 if event  $Z = z$  holds true; otherwise, it returns 0.

An SCM  $M$  is a tuple  $\langle \mathbf{V}, \mathbf{U}, \mathcal{F}, P(\mathbf{U}) \rangle$ , where  $\mathbf{V}$  is a set of endogenous variables and  $\mathbf{U}$  is a set of exogenous variables (Pearl, 2000; Bareinboim et al., 2022).  $\mathcal{F}$  is a set of functions s.t. each  $f_V \in \mathcal{F}$  decides values of an endogenous variable  $V \in \mathbf{V}$  taking as argument a combination of other variables in the system. That is,  $V \leftarrow f_V(\mathbf{PA}_V, \mathbf{U}_V)$ ,  $\mathbf{PA}_V \subseteq \mathbf{V}$ ,  $\mathbf{U}_V \subseteq \mathbf{U}$ . Values of exogenous variables  $U \in \mathbf{U}$  are drawn from the exogenous distribution  $P(\mathbf{U})$ . Naturally,  $M$  induces an *observational distribution*  $P(\mathbf{V})$ . An intervention on a subset  $\mathbf{X} \subseteq \mathbf{V}$ , denoted by  $\text{do}(\mathbf{x})$ , is an operation where values of  $\mathbf{X}$  are set to constants  $\mathbf{x}$ , replacing the functions  $\{f_X : \forall X \in \mathbf{X}\}$  that would normally determine their values. For an SCM  $M$ , let  $M_{\mathbf{x}}$  be a submodel of  $M$  induced by intervention  $\text{do}(\mathbf{x})$ .

For a set  $\mathbf{Y} \subseteq \mathbf{V}$ , the *interventional distribution*  $P_{\mathbf{x}}(\mathbf{Y})$  induced by  $\text{do}(\mathbf{x})$  is defined as the joint distribution over  $\mathbf{Y}$  in the submodel  $M_{\mathbf{x}}$ , i.e.,  $P_{\mathbf{x}}(\mathbf{Y}; M) \triangleq P(\mathbf{Y}; M_{\mathbf{x}})$ .

## 2 CHALLENGES OF CAUSAL INCONSISTENCY

We will focus on the policy evaluation problem of an agent operating in a Markov Decision Process (MDP) (Puterman, 1994) over a series of interventions  $t = 1, 2, \dots$ . For every time step  $t$ , the agent observes the current state  $S_t$ , performs an action  $\text{do}(X_t)$ , receives a subsequent reward  $Y_t$ , and moves to the next state  $S_{t+1}$ . Values of the action  $X_t$  are selected by sampling from a stationary policy  $\pi(x | s)$ , which is a function mapping from the domain of the observed state  $S_t$  to the probability space over the domain of action  $X_t$ . Let  $\mathbf{U}_t$  be an unobserved noise independently drawn from an exogenous distribution  $P(\mathbf{U})$ . Values of the reward  $Y_t$  and the next state  $S_{t+1}$  are, respectively, determined by structural functions  $y_t \leftarrow f_Y(s_t, x_t, \mathbf{u}_t)$  and  $s_{t+1} \leftarrow f_S(s_t, x_t, \mathbf{u}_t)$ , taking as input the current state  $S_t$ , action  $X_t$ , and latent noise  $\mathbf{U}_t$ ; values of  $S_1$  are drawn from an initial distribution  $P(S_1)$ . We will consistently use  $\mathcal{X}$ ,  $\mathcal{S}$ , and  $\mathcal{Y}$  to denote the domain of every action  $X_t$ , state  $S_t$ , and reward  $Y_t$ . Like a standard discrete MDP, domains of actions  $\mathcal{X}$  and states  $\mathcal{S}$  are assumed to be finite; rewards are bounded in a real interval  $\mathcal{Y} \triangleq [a, b] \subset \mathbb{R}$ . Naturally, the agent operating in this environment defines an interventional distribution  $P_{\pi}$  summarizing the consequences of its actions.

Fig. 1a shows a graphical representation (for now, without the highlighted bi-directed arrows) of this data-generating process where nodes represent observed variables and directed arrows represent the functional relationships between them. For every time step  $t > 1$ , the current state  $S_t$  “block” all pathways from previous nodes (e.g.,  $S_{t-1}$ ) to the future nodes (e.g.,  $S_{t+1}$ ) (Pearl, 2000, Def. 1.2.3). Applying the d-separation rules leads to the following independence relationships in distribution  $P_{\pi}$ .

**Definition 1** (Markov Property (Puterman, 1994)). For a distribution  $P_*$  over a sequence of states  $S_1, S_2, \dots$ , actions  $X_1, X_2, \dots$ , and rewards  $Y_1, Y_2, \dots$ , the Markov property holds if for every  $t = 1, 2, \dots$ ,  $(\bar{S}_{1:t-1}, \bar{X}_{1:t-1}, \bar{Y}_{1:t-1} \perp\!\!\!\perp \bar{X}_{t:\infty}, \bar{S}_{t+1:\infty}, \bar{Y}_{t:\infty} | S_t)$  with regard to distribution  $P_*$ .

It follows from Def. 1 that for any horizon  $T$ , the distribution generated by a policy  $\pi$  factories as

$$P_{\pi}(\bar{\mathbf{x}}_{1:T}, \bar{\mathbf{s}}_{1:T}, \bar{\mathbf{y}}_{1:T}) = P(s_1) \prod_{t=1}^T \pi(x_t | s_t) \mathcal{T}(s_t, x_t, s_{t+1}) \mathcal{R}(s_t, x_t, y_t) \quad (1)$$

where the transition distribution  $\mathcal{T}$  and the reward distribution  $\mathcal{R}$  are interventional queries given by

$$\mathcal{T}(s_t, x_t, s_{t+1}) = P_{x_t}(s_{t+1} | s_t) = \int_{\mathbf{u}_t} \mathbb{1}_{s_{t+1}=f_S(s_t, x_t, \mathbf{u}_t)} P(\mathbf{u}_t) \quad (2)$$

$$\mathcal{R}(s_t, x_t, y_t) = P_{x_t}(y_t | s_t) = \int_{\mathbf{u}_t} \mathbb{1}_{y_t=f_Y(s_t, x_t, \mathbf{u}_t)} P(\mathbf{u}_t) \quad (3)$$

For convenience, we write the reward function  $\mathcal{R}(s, x)$  as the expected value  $\sum_y y \mathcal{R}(s, x, y)$ . Fix a discounted factor  $\gamma \in [0, 1]$ . A common objective for an agent is to optimize its cumulative return  $R_t = \sum_{i=0}^{\infty} \gamma^i Y_{t+i}$ . In analysis, we often evaluate the state value function  $V_{\pi}(s)$ , which is the expected return given the agent’s starting state  $S_t = s$ . That is,  $V_{\pi}(s) = \mathbb{E}_{\pi}[R_t | S_t = s]$ . A similar state-action value function  $Q_{\pi}(s, x)$  is defined as the expected return starting from state  $s$ , taking action  $x$  and thereafter following policy  $\pi$ , i.e.,  $Q_{\pi}(s, x) = \mathbb{E}_{X_t \leftarrow x, \pi}[R_t | S_t = s]$ . One could recursively evaluate the value function of any state  $s$  using the *Bellman Equation* (Bellman, 1966):

$$V_{\pi}(s) = \sum_x \pi(x | s) \left( \mathcal{R}(s, x) + \gamma \sum_{s'} \mathcal{T}(s, x, s') V_{\pi}(s') \right) \quad (4)$$

Similarly, an analogous equation for the state-action value function is

$$Q_{\pi}(s, x) = \mathcal{R}(s, x) + \gamma \sum_{s'} \mathcal{T}(s, x, s') V_{\pi}(s') \quad (5)$$

**Off-Policy Evaluation** Despite the effectiveness of planning algorithms, they require detailed parametrization of the transition distribution  $\mathcal{T}$  and the reward function  $\mathcal{R}$ , which are not accessible in many real-world applications. This means that a learning process must take place. A common

162 approach is off-policy learning, where the agent has access to observed trajectories generated by a  
 163 *behavioral policy*  $f_X$ , different from the target policy  $\pi$ , operating in the same environment. More  
 164 specifically, for every time step  $t$ , the behavioral policy selects an action  $X_t \leftarrow f_X(s_t, \mathbf{u}_t)$  based on  
 165 the current state  $S_t = s_t$  and latent noise  $\mathbf{U}_t = \mathbf{u}_t$ . Fig. 1a shows the graphical representation of  
 166 the data-generating process of the behavior policy; the added bi-directed arrows, e.g.,  $X_t \leftrightarrow Y_t$ ,  
 167 indicate the presence of an unobserved confounder  $U \in \mathbf{U}_t$  affecting both the action  $X_t$  and outcome  
 168  $Y_t$ . We summarize observed trajectories of the behavior policy using the observational distribution  $P$ .

169 Off-policy evaluation attempts to estimate the effects of a candidate policy  $\pi(x|s)$  from the observa-  
 170 tional data generated by the behavior policy  $f_X$ . Standard off-policy methods focus on the identifiable  
 171 setting where the target transition distribution  $\mathcal{T}$  and reward function  $\mathcal{R}$  remain consistent in both the  
 172 interventional  $P_\pi$  and observational distribution  $P$ . Formally,

173 **Definition 2** (Causal Consistency). For an interventional distribution  $P_\pi$  and an observational  
 174 distribution  $P$  satisfying the Markov property (Def. 1), the Causal Consistency holds with regard to  
 175  $P_\pi$  and  $P$  if the following statement holds, for every time step  $t = 1, 2, \dots$ ,

$$176 P_{x_t}(s_{t+1} | s_t) = P(s_{t+1} | s_t, x_t), \quad \text{and} \quad P_{y_t}(y_t | s_t) = P(y_t | s_t, x_t) \quad (6)$$

178 When Def. 2 holds, the learner could recover the parametrization of the transition distribution  $\mathcal{T}$   
 179 and reward function  $\mathcal{R}$  from the observational data, following the identification formula in Eq. (6).  
 180 Several off-policy algorithms have been proposed to estimate the effect of candidate policies from  
 181 finite observations under causal consistency (Watkins, 1989; Watkins & Dayan, 1992; Swaminathan  
 182 & Joachims, 2015; Jiang & Li, 2016; Precup et al., 2000; Munos et al., 2016).

183 There exist graphical criteria in the literature (Pearl & Robins, 1995; Shpitser et al., 2010; Perković  
 184 et al., 2015) to evaluate whether causal consistency (Def. 2) holds from causal knowledge of the  
 185 environment, including the celebrated *backdoor* criterion (Pearl, 2000, Def. 3.3.1). However, in  
 186 many practical applications, causal consistency could be fragile and does not necessarily hold due to  
 187 some violations in the generative process. These include: (1) there exists an unobserved confounder  
 188 affecting the action  $X_t$  and subsequent outcomes  $Y_t, S_{t+1}$  simultaneously (blue, dashed arrows in  
 189 Fig. 1a); (2) there is no overlap in the support between the target and behavior policies, i.e., the  
 190 propensity score  $P(x_t | s_t) = 0$  for some state-action pair  $s_t, x_t$ . When either of these violations  
 191 occurs, applying standard off-policy methods may fail to recover the expected return of the target  
 192 policy, leading to estimation bias. The following example illustrates such challenges.

193 **Example 1** (Windy Gridworld). Consider a Windy Gridworld described in Fig. 1b, where the red dot  
 194 represents the agent and the green square represents the goal state. The agent can take five actions  
 195  $X_t$  - up, down, right, left, and stay-put. However, the agent’s movement is affected by the  
 196 wind; the direction of the wind  $U_t$  includes - north, south, west, east, and no-wind. For  
 197 every time step, the agent receives a constant reward  $Y_t \leftarrow -1$ . The next state of the agent is shifted  
 198 by both its action and the wind direction through the mechanism  $S_{t+1} \leftarrow S_t + X_t + U_t$ .

199 Our goal is to evaluate the expected return of a target policy  $\pi^*$  described in Fig. 2a, which consistently  
 200 moves towards the goal state regardless of the wind direction. As an input, we have access to observed  
 201 trajectories generated by a behavior policy  $X_t \leftarrow f_X(S_t, U_t)$ , which could sense the wind and select  
 202 an action accordingly. For example, when the agent is located in the top-left corner ( $S_t = (0, 0)$ )  
 203 and the wind is blowing south ( $U_t = (0, 1)$ ), the behavior policy will decide to move right  
 204 ( $X_t = (1, 0)$ ) so that the agent could get close to the center ( $S_{t+1} = (1, 1)$ ).

205 Figs. 2b to 2d shows the value function estimation obtained by standard off-policy methods, including  
 206 Q-Learning, one-step Temporal Difference (TD), and Eligibility Traces (TD( $\lambda$ )). We also include  
 207 in Fig. 2e the ground truth value function computed from the underlying model parameters. The  
 208 simulation reveals that standard off-policy evaluation deviates from the ground truth return. In this  
 209 observational data, the wind direction  $U_t$  is thus an unobserved confounder affecting both the action  
 210  $X_t$  and next state  $S_{t+1}$ , violating causal consistency. See Appendix C for additional discussions.

## 212 2.1 PARTIAL CAUSAL IDENTIFICATION IN MDPs

213 For the remainder of this section, we will introduce partial identification methods for off-policy  
 214 evaluation that is robust to the unobserved confounding and no-overlap. For every time step  $t =$   
 215  $1, 2, \dots$ , let the reward  $Y_t$  be bounded in a real interval  $[a, b]$ . By applying a similar bounding strategy

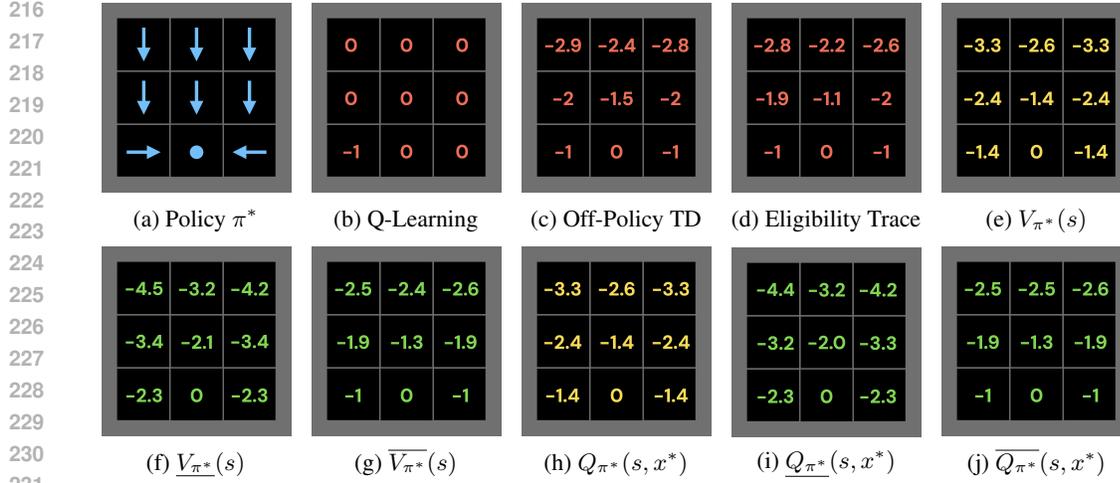


Figure 2: (a) The target policy  $\pi^*$  selecting an action based on the agent’s location. (b - d) Value function estimation was obtained by standard off-policy methods. (e - g) The ground-truth state value function computed from the model parametrization and its lower and upper bounds estimated using the extended Bellman equation in Thm. 1. (h - j) The ground-truth state-action value function computed from the model parametrization for actions  $x^* \leftarrow \pi^*(s)$  selected by the target policy and its lower and upper bounds computed from the extended Bellman equation in Thm. 2

in (Manski, 1990; Zhang & Bareinboim, 2019; Joshi et al., 2024), we derive the following bounds over the transition distribution  $\mathcal{T}$  and reward function  $\mathcal{R}$ , for every realization  $(s, x, s') \in \mathcal{S} \times \mathcal{X} \times \mathcal{S}$ ,

$$\mathcal{T}(s, x, s') \in \left[ \tilde{\mathcal{T}}(s, x, s') P(x | s), \tilde{\mathcal{T}}(s, x, s') P(x | s) + P(\neg x | s) \right] \quad (7)$$

$$\mathcal{R}(s, x) \in \left[ \tilde{\mathcal{R}}(s, x) P(x | s) + aP(\neg x | s), \tilde{\mathcal{R}}(s, x) P(x | s) + bP(\neg x | s) \right] \quad (8)$$

Among the above quantities,  $P(x | s)$  stands for the propensity score  $P(X_t = x | S_t = s)$  and  $P(\neg x | s) = 1 - P(x | s)$ ;  $\tilde{\mathcal{T}}$  and  $\tilde{\mathcal{R}}$  are the nominal transition distribution and reward function computed from the observational distribution as follows:

$$\tilde{\mathcal{T}}(s, x, s') = P(S_{t+1} = s' | S_t = s, X_t = x), \quad \tilde{\mathcal{R}}(s, x) = \mathbb{E}[Y_t | S_t = s, X_t = x] \quad (9)$$

In order to bound the value function  $V_\pi(s)$  at state  $s$  induced by a candidate policy  $\pi$ , one could minimize/maximize the optimization program using the Bellman’s equation in Eq. (4) as the objective function, subject to constraints in Eqs. (7) and (8). Interestingly, this optimization problem is equivalent to a linear program; solving it leads to the following *extended Bellman equation*.

**Theorem 1** (Causal Bellman Equation). *For an MDP environment  $M$  with reward  $Y_t \in [a, b] \subseteq \mathbb{R}$ , for any policy  $\pi(x | s)$ , its state value function  $V_\pi(s) \in [V_\pi(s), \overline{V}_\pi(s)]$  for every state  $s \in \mathcal{S}$ , where bounds  $\underline{V}_\pi, \overline{V}_\pi$  are solutions given by the following dynamic programs,<sup>1</sup>*

$$\langle \underline{V}_\pi(s), \overline{V}_\pi(s) \rangle = \sum_x P(x | s) \left( \pi(x | s) \left( \tilde{\mathcal{R}}(s, x) + \gamma \sum_{s', x'} \tilde{\mathcal{T}}(s, x, s') \langle \underline{V}_\pi(s'), \overline{V}_\pi(s') \rangle \right) \right) \quad (10)$$

$$+ \pi(\neg x | s) \left( \langle a, b \rangle + \gamma \langle \min_{s'} \underline{V}_\pi(s'), \max_{s'} \overline{V}_\pi(s') \rangle \right) \quad (11)$$

Thm. 1 can be seen as an extension of the Bellman equation using the confounded observational distribution with no-overlap. For instance, in the lower bound  $\underline{V}_\pi(s)$ , Eq. (10) follows the standard iterative step in Bellman equation in Eq. (4), measuring the expected return when the target policy’s

<sup>1</sup> $\langle a, b \rangle$  is a vector containing a lower bound  $a$  and an upper bound  $b$ . We highlight quantities that are different from the standard Bellman Equation.

action coincides with the observed action selected by the behavior policy; Eq. (11) could be thought as a regularizing term measuring the uncertainty due to unobserved confounding. Finally, both terms are weighted by the nominal propensity score  $P(x | s) = P(X_t = x | S_t = s)$ . The same derivation also applies to the upper bound  $\overline{V}_\pi(s)$ . An analogous extended Bellman equation bounding the state-action value function from the observational distribution can also be derived as follows.

**Theorem 2** (Causal Bellman Equation). *For an MDP environment  $M$  with reward signals  $Y_t \in [a, b] \subseteq \mathbb{R}$ , for any policy  $\pi(x | s)$ , its state-action value function  $Q_\pi \in [Q_\pi(s, x), \overline{Q}_\pi(s, x)]$  for any state-action pair  $(s, x) \in \mathcal{S} \times \mathcal{X}$ , where bounds  $\underline{Q}_\pi, \overline{Q}_\pi$  are given by as follows,*

$$\left\langle \underline{Q}_\pi(s, x), \overline{Q}_\pi(s, x) \right\rangle = P(x | s) \left( \tilde{\mathcal{R}}(s, x) + \gamma \sum_{s', x'} \tilde{\mathcal{T}}(s, x, s') \left\langle \underline{V}_\pi(s'), \overline{V}_\pi(s') \right\rangle \right) \quad (12)$$

$$+ P(-x | s) \left( \langle a, b \rangle + \gamma \left\langle \min_{s'} \underline{V}_\pi(s'), \max_{s'} \overline{V}_\pi(s') \right\rangle \right) \quad (13)$$

Among the bounds in Thm. 2, Eq. (12) is the standard iterative step of the Bellman equation in Eq. (5), weighted by the score  $P(x | s)$ . It estimates the expected return of performing action  $\text{do}(x)$  at state  $s$  when such action matches the one selected by the behavior policy. Eq. (13) is a regularized term accounting for uncertainties when the intervention  $\text{do}(x)$  is not observed in the offline data. Since Thms. 1 and 2 are closed-form solutions of optimization programs and the observational constraints in Eqs. (7) and (8) are tight, the extended Bellman’s equation bounds are sharp from offline data and Markov property. This means they cannot be improved without additional assumptions.

**Example 2** (Windy Gridworld Continued). Consider again the Windy Gridworld described in Example 1. We compute the lower and upper bounds over the state value function following the extended Bellman equation in Thm. 1, and provide them in Figs. 2f to 2g. We also include in Fig. 2h the ground truth state-action value function for the action  $x^* \leftarrow \pi^*(s)$  selected by the target policy. The corresponding lower and upper bounds are shown in Figs. 2i to 2j, following the algorithmic procedure described in Thm. 2. The analysis reveals that the derived bounds are consistent with the ground truth value functions, corroborating the sufficiency of our proposed approach.

### 3 CONFOUNDING ROBUST ELIGIBILITY TRACES

The extended Bellman equations described so far require one to have precise estimations for the full models of the nominal transition distribution  $\mathcal{T}_{\text{obs}}$ , reward function  $\mathcal{R}_{\text{obs}}$ , and the propensity score  $P(x | s)$ . This section will introduce novel model-free algorithms, using eligibility traces (Sutton, 1988), to bound value functions from finite observational samples.

We consider the episodic framework, where the agent interacts with the environment for repeated episodes  $n = 1, 2, 3, \dots$ ; each episode contains a finite number of time steps  $t = 1, 2, \dots, T_n$ . At each episode, the environment starts at state  $s_1$  following the initial distribution  $P(S_1)$ . At each time step  $t$ , taking the observed state  $s_t$  of the environment as input, the behavior policy selects an action  $x_t$ . In response to intervention  $\text{do}(x_t)$ , the environment produces a subsequent reward  $y_t$  and moves to the next observed state  $s_{t+1}$ . If the next state  $s_{t+1}$  is *terminal*, the episode terminates at time step  $T_n = t + 1$ ; the learner receives observational data  $\{\bar{x}_{1:T_n-1}, \bar{s}_{1:T_n}, \bar{y}_{1:T_n-1}\}$ .

#### 3.1 CAUSAL TEMPORAL DIFFERENCE

We first introduce a novel augmentation procedure on the celebrated temporal difference (TD, (Sutton, 1988; Precup et al., 2000)) that allows one to estimate the bounds over state value functions, which we call the *causal temporal difference* (C-TD). Fig. 3 shows the backup diagram illustrating the idea of our proposed algorithm. Similar to the standard off-policy TD, our algorithm will update the estimation of state value functions  $\underline{V}_\pi, \overline{V}_\pi$  using the sampled trajectories of transitions in the observational data. It could use a finite number of  $n$ -step trajectories or the entire trajectory. Different from the standard off-policy TD, our proposed algorithm does not weight each step of the transition using importance sampling (or equivalently, inverse propensity weighting) since the true behavior policy  $f_X$  (propensity

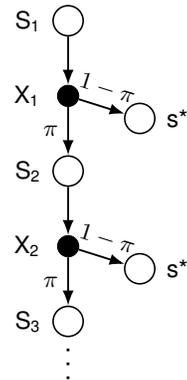


Figure 3: Backup diagram for C-TD ( $\lambda$ ).

**Algorithm 1** Causal Temporal Difference (C-TD ( $\lambda$ ))**Require:** Observational data  $\mathcal{D}$  and a candidate policy  $\pi(x | s)$ .

- 1: Update the eligibility traces for all state
- $s$
- ,

$$e_t(s) = \begin{cases} \gamma\lambda\pi(x_{t-1} | s_{t-1})e_{t-1}(s) & \text{if } s \neq s_t \\ \gamma\lambda\pi(x_{t-1} | s_{t-1})e_{t-1}(s) + 1 & \text{if } s = s_t \end{cases} \quad (15)$$

where  $\lambda \in [0, 1]$  is an eligibility trace decay factor.

- 2: Compute the temporal difference error

$$\delta_t = \pi(x_t | s_t)(y_t + \gamma V_t(s_{t+1})) + \pi(\neg x_t | s_t)(w + \gamma V_t(s^*)) - V_t(s_t) \quad (16)$$

- 3: Update the value function
- $V_{t+1}(s) \leftarrow V_t(s) + \alpha e_t(s)\delta_t$
- for all state
- $s$
- .

**Algorithm 2** Causal Tree-Backup (C-TB ( $\lambda$ ))**Require:** Observational data  $\mathcal{D}$  and a candidate policy  $\pi(x|s)$ .

- 1: Update the eligibility traces for all state-action pairs
- $s, x$
- ,

$$e_t(s, x) = \begin{cases} \gamma\lambda\pi(x_t | s_t)\mathbb{1}_{x_{t-1}=x}e_{t-1}(s, x) & \text{if } s \neq s_t \\ \gamma\lambda\pi(x_t | s_t)\mathbb{1}_{x_{t-1}=x}e_{t-1}(s, x) + 1 & \text{if } s = s_t \end{cases} \quad (17)$$

where  $\lambda \in [0, 1]$  is an eligibility trace decay factor.

- 2: Compute the temporal difference error for every action
- $x$

$$\delta_t(x) = \begin{cases} y_t + \gamma \sum_{x'} \pi(x | s_{t+1})Q_t(s_{t+1}, x') - Q_t(s_t, x) & \text{if } x = x_t \\ w + \gamma \sum_{x'} \pi(x' | s^*)Q_t(s^*, x') - Q_t(s_t, x) & \text{if } x \neq x_t \end{cases} \quad (18)$$

- 3: Update the action-value function
- $Q_{t+1}(s, x) \leftarrow Q_t(s, x) + \alpha e_t(s, x)\delta_t(x)$
- for all
- $s, x$
- .

score) is not recoverable from the observational data. Instead, C-TD weights each transition using the target policy  $\pi$  and adjusts for the misalignment between the target and behavior policies using an overestimation/underestimation of value function at state  $s^*$ . Such  $s^*$  is set as the best-case state associated with the highest value in our current estimation when computing upper bounds and the worst-case state estimate for lower bounds.

To formally introduce the estimation algorithm, we first introduce some necessary notations. Let  $\mathcal{N}(s)$  denote the set of indices of episodes containing a state  $s \in \mathcal{S}$ , and let  $\mathbf{t}_n(s)$  be the collection of time steps in the  $n$ -th episode such that for every  $t \in \mathbf{t}_n(s)$ ,  $s_t = s$ . For any time step  $t$ , let  $\pi_t = \pi(x_t | s_t)$  and  $\neg\pi_t = 1 - \pi(x_t | s_t)$ . We iteratively define the estimator for bounds over the state value function  $V_\pi(s)$  as follows, for any state  $s \in \mathcal{S}$ ,

$$\widehat{V}_\pi(s) = \frac{1}{N} \sum_{n \in \mathcal{N}(s)} \sum_{t \in \mathbf{t}_n(s)} \sum_{k=0}^{T_n-t} \gamma^k \left( \pi_{t+k} y_{t+k} + \neg\pi_{t+k} (w + \gamma V(s^*)) \right) \prod_{i=t}^{t+k-1} \pi_i, \quad (14)$$

Among the above equation,  $N$  represents the total number of occurrences for the even  $s_t = s$  in the observational data. we set parameters  $w = a$  and  $V(s^*) = \min_s V(s)$  when estimating the lower bound  $\underline{V}_\pi(s)$ ; parameters  $w = b$  and  $V(s^*) = \max_s V(s)$  for the upper bound  $\overline{V}_\pi(s)$ .

An eligibility-trace version of our proposed estimation strategy is described Alg. 1. The algorithm keeps track of eligibility traces for every state in a similar manner to standard off-policy temporal difference algorithms. The main difference is that here the eligibility trace is multiplied by the target policy  $\pi(x_{t-1} | s_{t-1})$  and a decay-rate  $\lambda$ , not including the nominal propensity score  $P(x_{t-1} | s_{t-1})$ . When computing the temporal difference error, the algorithm adjusts for the misalignment between the target and behavior policies by adding a regularized term  $w + \gamma V_t(s^*)$ , weighted by the probability  $1 - \pi(x_t | s_t)$ . We describe in Alg. 1 a version of C-TD ( $\lambda$ ) using *online update*. This means that the bounds estimate over value functions are updated at every time step. The *offline* version of the algorithm will use the same temporal difference error and eligibility traces. However, the update only occurs at the end of each episode; the increments and decrements are accumulated on the side, and the value function estimates do not change during the episode.

**Theorem 3.** For any behavior policy, for any choice of  $\lambda \in [0, 1]$  that does not depend on the actions chosen at each state, let parameters  $w$  and  $s^*$  be defined as follows: (1) Lower Bound  $\underline{V}_\pi$ :  $w = a$  and  $s^* = \arg \min_s V_t(s)$ ; (2) Upper Bound  $\overline{V}_\pi$ :  $w = b$  and  $s^* = \arg \max_s V_t(s)$ . Then, Alg. 1 with offline updating converges with probability 1 to lower bound  $\underline{V}_\pi$  and upper bound  $\overline{V}_\pi$ , respectively, under the usual step-size conditions on  $\alpha$ .

The proof of Thm. 3 first shows a contraction property for estimates  $\widehat{V}_\pi$ , and then follows the general convergence theorem in (Jaakkola et al., 1994).

### 3.2 CAUSAL TREE BACKUP

The algorithm described so far focuses on the estimation of the state value functions. We next introduce a novel algorithm to bound the state-action value function  $Q_\pi$  from finite observations.

Our algorithm is based on an augmentation on the standard tree backup (TB (Precup et al., 2000)), which we call the *causal tree backup* (C-TB ( $\lambda$ )). The main idea of this new algorithm is illustrated in the backup diagram of Fig. 4. Similar to the standard tree backup, our algorithm updates the value estimates for the action selected by the behavior policy at each time step based on the subsequent reward and the current estimation for the value of the next state. The algorithm then forms a new estimate for the target value function, using the old value estimates for the actions not observed in the observational data and the new estimated value for  $t$ -th action taken by the behavior policy. On the other hand, the main differences include the following. (1) Eligibility traces will not only be weighted by the target policy  $\pi(x_t | s_t)$  using the observed trajectories, but also an indicator function  $\mathbb{1}_{x_{t-1}=x}$  returning 1 if the previous action  $x_{t-1}$  coincides with the target action  $x$ . (2) When the behavior policy takes the same action  $x_t = x$  as the target action, the update follows standard TB and uses the next sampled state  $s_t$ ; when the sampled action  $x_t \neq x$  differs from the target, our algorithm updates, instead, using the value function associated with the next worst-case or best-case state  $s^*$ , corresponding to the estimation of the lower bound and upper bound respectively. The  $n$ -step causal tree-backup estimator is defined as

$$\widehat{Q}_\pi(s, x) = \frac{1}{N} \sum_{n \in \mathcal{N}(s)} \sum_{t \in \mathcal{t}(s)} \gamma^n Q(s_{t+n}, x_{t+n}) \prod_{i=t}^{t+n-1} \pi_{i+1} \mathbb{1}_{x_i=x} + \sum_{k=t}^{t+n} \gamma^{k-t+1} \prod_{i=t}^{t+k-1} \pi_{i+1} \mathbb{1}_{x_i=x} \cdot \left( \mathbb{1}_{x_k=x} \left( y_k + \sum_{x' \neq x} \pi(x' | s_{k+1}) Q(s_{k+1}, x') \right) + \mathbb{1}_{x_k \neq x} \left( w + \sum_{x'} \pi(x' | s^*) Q(s^*, x') \right) \right) \quad (19)$$

The above tree backup estimator also has a simple incremental implementation using eligibility traces. An online version of this implementation is shown in Fig. 4.

**Theorem 4.** For any behavior policy, for any choice of  $\lambda \in [0, 1]$  that does not depend on the actions chosen at each state, let parameters  $w$  and  $s^*$  be defined as follows: (1) Lower Bound  $\underline{Q}_\pi$ :  $w = a$  and  $s^* = \arg \min_s \sum_{x'} \pi(x' | s) Q_t(s, x')$ ; (2) Upper Bound  $\overline{Q}_\pi$ :  $w = b$  and  $s^* = \arg \max_s \sum_{x'} \pi(x' | s) Q_t(s, x')$ . Then, Alg. 2 with offline updating converges with probability 1 to lower bound  $\underline{Q}_\pi$  and upper bound  $\overline{Q}_\pi$ , respectively, under the usual step-size conditions on  $\alpha$ .

The proof of the above theorem relies on a contraction property on the estimates  $\widehat{Q}_\pi$  and follows from the general convergence theorem in (Jaakkola et al., 1994).

## 4 EXPERIMENTS

We demonstrate our algorithms using different behavior policies in the Windy Gridworld described in Example 1. Overall, we found that simulation results support our findings, and the proposed algorithms consistently obtain informative bounds over value functions. Experiment 1 evaluates the performance of our bounding strategy in the presence of unobserved confounding. Experiment 2 uses data collected from a deterministic sub-optimal policy, violating the overlap. All experiments use

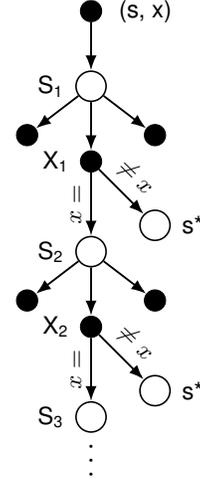
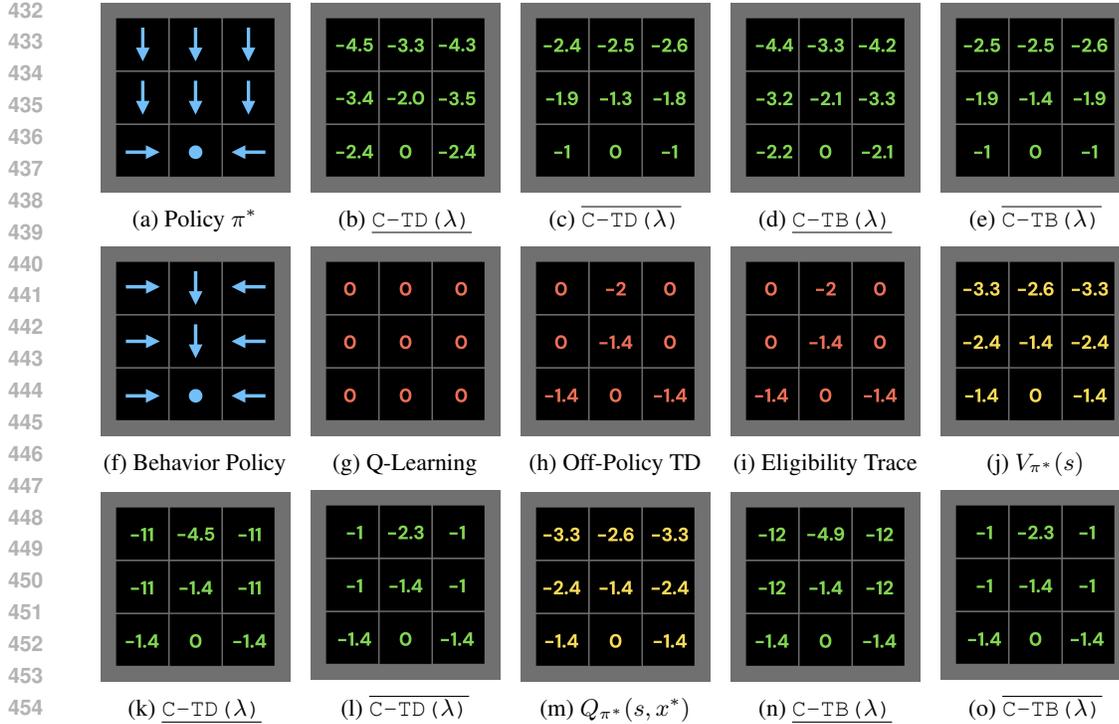


Figure 4: Backup diagram for C-TB ( $\lambda$ ).



456  
457  
458  
459  
460

Figure 5: Simulation results comparing causally enhanced off-policy algorithms using eligibility traces ( $\underline{C-TD}(\lambda)$  and  $\overline{C-TB}(\lambda)$ ) with standard off-policy methods. The offline data are collected from (a - e) a confounded behavior policy affected by the unobserved confounder; and (f - o) a deterministic behavior policy following sub-optimal actions.

461  
462  
463

$5 \times 10^4$  offline observational samples, meaning that error bars are not significant, hence, not explicitly shown; the decay factor  $\lambda = 0.5$ . See Appendix B for more details on the experimental setup.

464  
465  
466  
467  
468  
469  
470

**Experiment 1.** Consider again the learning setting described in Example 1 where the offline data is contaminated with unobserved confounding bias, and the behavior policy selects actions based on the agent’s state and the latent wind direction. We apply  $\underline{C-TD}(\lambda)$  to derive bounds over the state value function  $V_{\pi^*}(s)$  and provide them in Figs. 5b and 5c. We also compute the bounds over the state-action value function  $Q_{\pi^*}(s, x^*)$  for actions  $x^* \leftarrow \pi^*(s)$  using  $\underline{C-TB}(\lambda)$ ; the simulation results are shown in Figs. 5d and 5e. The analysis reveals that our algorithm consistently recovers the closed-form bounds containing the ground-truth value functions, as previously shown in Fig. 2.

471  
472  
473  
474  
475  
476  
477  
478  
479  
480  
481

**Experiment 2.** For the Windy Gridworld environment described in Example 1, suppose the data is now collected by a deterministic behavior policy that always first moves towards the center and then moves down toward the goal; its parametrization is provided in Fig. 5f. This means that the overlap does not hold when the agent is located on either side of the top half of the board. We apply standard off-policy algorithms to evaluate the effect of the target policy  $\pi^*$  of Fig. 5a and provide their evaluations in Figs. 5g to 5i. The propensity score is truncated using a small positive real  $0 < \epsilon < 1$  if  $P(x | s) = 0$ . We also compute bounds over the target value functions using our proposed algorithms,  $\underline{C-TD}(\lambda)$  and  $\overline{C-TB}(\lambda)$ , and provide their evaluations in Figs. 5k to 5l and Figs. 5n to 5o respectively. By comparing with the ground-truth values in Figs. 5j and 5m, we found that  $\underline{C-TD}(\lambda)$  and  $\overline{C-TB}(\lambda)$  can consistently obtain informative bounds; as expected, standard off-policy methods are not robust against no-overlap and deviate significantly from the target effects.

## 482 483 484 485

## 5 CONCLUSION

This paper investigates off-policy evaluation in Markov Decision Processes from offline data collected by a different behavior policy, where unobserved confounding bias and no-overlap cannot be ruled

486 out *a priori*. This leads to violations of causal consistency (Def. 2), which could pose significant  
 487 challenges to standard off-policy algorithms. We first extend the celebrated Bellman’s equation  
 488 to derive informative bounds over value functions from the observational data, which are robust  
 489 against bias due to the presence of unobserved confounding and no-overlap. Based on these extended  
 490 equations, we propose two novel model-free off-policy algorithms using eligibility traces – one based  
 491 on the standard temporal difference (C-TD ( $\lambda$ )), and the other based on the tree-backup (C-TB ( $\lambda$ )).  
 492 These algorithms permit us to bound value functions from finite observations consistently.

## 493 494 ETHICS STATEMENT

496 This paper investigates the theoretical framework of robust off-policy evaluation from biased offline  
 497 data generated by a different behavior. Since unobserved confounding or no-overlap cannot be  
 498 ruled out *a priori*, the agent’s system dynamics in the environment cannot be fully identified from  
 499 the offline data. To address this challenge, we proposed novel off-policy algorithms that allow the  
 500 agent to derive informative bounds over value functions induced by a target policy from biased  
 501 offline data. A positive impact of this work is that we address the potential risk of policy learning  
 502 from offline data with the presence of unobserved confounding. Our framework is inherently robust  
 503 against confounding bias and may apply to various consequential domains involving complex human  
 504 interactions, including healthcare, marketing, finance, and autonomous driving. More broadly,  
 505 automated decision systems using causal inference methods prioritize safety and robustness during  
 506 their learning processes. Such requirements are increasingly essential since black-box AI systems are  
 507 prevalent, and our understanding of their potential implications is still limited.

## 508 509 REPRODUCIBILITY STATEMENT

510 The complete proof of all theoretical results presented in this paper, including Thms. 1 to 4, is  
 511 provided in Appendix B. Detailed descriptions of the experimental setup are included in Appendix C.  
 512 Readers can find all appendices as part of the supplementary text after the “References” section. All  
 513 the experiments are synthetic and do not introduce any new assets. Windy Gridworld is implemented  
 514 based on the Gymnasium framework (Towers et al., 2024).

## 516 517 REFERENCES

- 518 Kishan Panaganti Badrinath and Dileep Kalathil. Robust reinforcement learning using least squares  
 519 policy iteration with provable performance guarantees. In *International Conference on Machine*  
 520 *Learning*, pp. 511–520. PMLR, 2021.
- 521 Patrick Bajari, C Lanier Benkard, and Jonathan Levin. Estimating dynamic models of imperfect  
 522 competition. *Econometrica*, 75(5):1331–1370, 2007.
- 523 A. Balke and J. Pearl. Bounds on treatment effects from studies with imperfect compliance. *Journal*  
 524 *of the American Statistical Association*, 92(439):1172–1176, September 1997.
- 525 Elias Bareinboim, Juan D Correa, Duligur Ibeling, and Thomas Icard. On pearl’s hierarchy and the  
 526 foundations of causal inference. In *Probabilistic and causal inference: the works of judea pearl*,  
 527 pp. 507–556. 2022.
- 528 Richard Bellman. Dynamic programming. *Science*, 153(3731):34–37, 1966.
- 529 Steven T Berry and Giovanni Compiani. An instrumental variable approach to dynamic models. *The*  
 530 *Review of Economic Studies*, 90(4):1724–1758, 2023.
- 531 Greg Brockman, Vicki Cheung, Ludwig Pettersson, Jonas Schneider, John Schulman, Jie Tang, and  
 532 Wojciech Zaremba. Openai gym. *arXiv preprint arXiv:1606.01540*, 2016.
- 533 David Bruns-Smith and Angela Zhou. Robust fitted-q-evaluation and iteration under sequentially  
 534 exogenous unobserved confounders. *arXiv preprint arXiv:2302.00662*, 2023.
- 535 Federico A Bugni. Bootstrap inference in partially identified models defined by moment inequalities:  
 536 Coverage of the identified set. *Econometrica*, 78(2):735–753, 2010.

- 540 Carlos Cinelli, Daniel Kumor, Bryant Chen, Judea Pearl, and Elias Bareinboim. Sensitivity analysis of  
541 linear structural causal models. In *International Conference on Machine Learning*, pp. 1252–1261,  
542 2019.
- 543 Michael J Dickstein and Eduardo Morales. What do exporters know? *The Quarterly Journal of*  
544 *Economics*, 133(4):1753–1801, 2018.
- 545 Mohammad Ghavamzadeh, Marek Petrik, and Yinlam Chow. Safe policy improvement by minimizing  
546 robust baseline regret. *Advances in Neural Information Processing Systems*, 29, 2016.
- 547 G Imbens and J Angrist. Estimation and identification of local average treatment effects. *Economet-*  
548 *rica*, 62:467–475, 1994.
- 549 Guido W Imbens and Donald B Rubin. Bayesian inference for causal effects in randomized experi-  
550 ments with noncompliance. *The annals of statistics*, pp. 305–327, 1997.
- 551 Garud N Iyengar. Robust dynamic programming. *Mathematics of Operations Research*, 30(2):  
552 257–280, 2005.
- 553 Tommi Jaakkola, Satinder Singh, and Michael Jordan. Reinforcement learning algorithm for partially  
554 observable markov decision problems. *Advances in neural information processing systems*, 7,  
555 1994.
- 556 Andrew Jesson, Alyson Douglas, Peter Manshausen, Nicolai Meinshausen, Philip Stier, Yarin  
557 Gal, and Uri Shalit. Scalable sensitivity and uncertainty analysis for causal-effect estimates of  
558 continuous-valued interventions. *arXiv preprint arXiv:2204.10022*, 2022.
- 559 Nan Jiang and Lihong Li. Doubly robust off-policy value evaluation for reinforcement learning. In  
560 Maria Florina Balcan and Kilian Q. Weinberger (eds.), *Proceedings of The 33rd International*  
561 *Conference on Machine Learning*, volume 48 of *Proceedings of Machine Learning Research*, pp.  
562 652–661, New York, New York, USA, 20–22 Jun 2016. PMLR. URL <http://proceedings.mlr.press/v48/jiang16.html>.
- 563 Shalmali Joshi, Junzhe Zhang, and Elias Bareinboim. Towards safe policy learning under partial  
564 identifiability: A causal approach. In *Proceedings of the AAAI Conference on Artificial Intelligence*,  
565 volume 38, pp. 13004–13012, 2024.
- 566 Nathan Kallus and Angela Zhou. Confounding-robust policy improvement. In *Proceedings of the*  
567 *32nd International Conference on Neural Information Processing Systems*, pp. 9289–9299, 2018.
- 568 Nathan Kallus and Angela Zhou. Confounding-robust policy evaluation in infinite-horizon reinforce-  
569 ment learning. In H. Larochelle, M. Ranzato, R. Hadsell, M. F. Balcan, and H. Lin (eds.), *Advances*  
570 *in Neural Information Processing Systems*, volume 33, pp. 22293–22304. Curran Associates, Inc.,  
571 2020.
- 572 Nathan Kallus, Xiaojie Mao, and Angela Zhou. Interval estimation of individual-level causal effects  
573 under unobserved confounding. In *The 22nd International Conference on Artificial Intelligence*  
574 *and Statistics*, pp. 2281–2290. PMLR, 2019.
- 575 Chinmaya Kausik, Yangyi Lu, Kevin Tan, Maggie Makar, Yixin Wang, and Ambuj Tewari. Offline  
576 policy evaluation and optimization under confounding. In *International Conference on Artificial*  
577 *Intelligence and Statistics*, pp. 1459–1467. PMLR, 2024.
- 578 Samir Khan, Martin Saveski, and Johan Ugander. Off-policy evaluation beyond overlap: partial  
579 identification through smoothness. *arXiv preprint arXiv:2305.11812*, 2023.
- 580 Shiau Hong Lim, Huan Xu, and Shie Mannor. Reinforcement learning in robust markov decision  
581 processes. *Advances in Neural Information Processing Systems*, 26, 2013.
- 582 Shie Mannor, Ofir Mebel, and Huan Xu. Robust mdps with k-rectangular uncertainty. *Mathematics*  
583 *of Operations Research*, 41(4):1484–1509, 2016.
- 584 C.F. Manski. Nonparametric bounds on treatment effects. *American Economic Review, Papers and*  
585 *Proceedings*, 80:319–323, 1990.

- 594 Hyungsik Roger Moon and Frank Schorfheide. Bayesian and frequentist inference in partially  
595 identified models. *Econometrica*, 80(2):755–782, 2012.
- 596
- 597 Eduardo Morales, Gloria Sheu, and Andrés Zahler. Extended gravity. *The Review of economic*  
598 *studies*, 86(6):2668–2712, 2019.
- 599 Rémi Munos, Tom Stepleton, Anna Harutyunyan, and Marc Bellemare. Safe and efficient off-policy  
600 reinforcement learning. In *Advances in Neural Information Processing Systems*, pp. 1054–1062,  
601 2016.
- 602
- 603 Hongseok Namkoong, Ramtin Keramati, Steve Yadlowsky, and Emma Brunskill. Off-policy policy  
604 evaluation for sequential decisions under unobserved confounding. *Advances in Neural Information*  
605 *Processing Systems*, 33:18819–18831, 2020.
- 606
- 607 Arnab Nilim and Laurent El Ghaoui. Robust control of markov decision processes with uncertain  
608 transition matrices. *Operations Research*, 53(5):780–798, 2005.
- 609 Andriy Norets and Xun Tang. Semiparametric inference in dynamic binary choice models. *Review of*  
610 *Economic Studies*, 81(3):1229–1262, 2014.
- 611
- 612 Kishan Panaganti, Zaiyan Xu, Dileep Kalathil, and Mohammad Ghavamzadeh. Robust reinforcement  
613 learning using offline data. *Advances in neural information processing systems*, 35:32211–32224,  
614 2022.
- 615 J. Pearl and J.M. Robins. Probabilistic evaluation of sequential plans from causal models with  
616 hidden variables. In P. Besnard and S. Hanks (eds.), *Proceedings of the Eleventh Conference on*  
617 *Uncertainty in Artificial Intelligence (UAI 1995)*, pp. 444–453. Morgan Kaufmann, San Francisco,  
618 1995.
- 619 Judea Pearl. *Causality: Models, Reasoning, and Inference*. Cambridge University Press, New York,  
620 2000.
- 621
- 622 Emilija Perković, Johannes Textor, Markus Kalisch, and Marloes H Maathuis. A complete generalized  
623 adjustment criterion. *arXiv preprint arXiv:1507.01524*, 2015.
- 624
- 625 Marek Petrik and Reazul Hasan Russel. Beyond confidence regions: Tight bayesian ambiguity sets  
626 for robust mdps. *Advances in neural information processing systems*, 32, 2019.
- 627 Lerrel Pinto, James Davidson, Rahul Sukthankar, and Abhinav Gupta. Robust adversarial rein-  
628 forcement learning. In *International conference on machine learning*, pp. 2817–2826. PMLR,  
629 2017.
- 630
- 631 Dale J Poirier. Revising beliefs in nonidentified models. *Econometric theory*, 14(4):483–509, 1998.
- 632
- 633 Doina Precup, Richard S. Sutton, and Satinder P. Singh. Eligibility traces for off-policy policy  
634 evaluation. In *Proceedings of the Seventeenth International Conference on Machine Learning*, pp.  
635 759–766, 2000.
- 636
- 637 Martin L Puterman. *Markov Decision Processes: Discrete Stochastic Dynamic Programming*. John  
Wiley & Sons, Inc., 1994.
- 638
- 639 Amy Richardson, Michael G Hudgens, Peter B Gilbert, and Jason P Fine. Nonparametric bounds  
640 and sensitivity analysis of treatment effects. *Statistical science: a review journal of the Institute of*  
641 *Mathematical Statistics*, 29(4):596, 2014.
- 642
- 643 Joseph P Romano and Azeem M Shaikh. Inference for identifiable parameters in partially identified  
econometric models. *Journal of Statistical Planning and Inference*, 138(9):2786–2807, 2008.
- 644
- 645 Paul R Rosenbaum. Sensitivity analysis in observational studies. *Encyclopedia of statistics in*  
646 *behavioral science*, 2005.
- 647
- Aurko Roy, Huan Xu, and Sebastian Pokutta. Reinforcement learning under model mismatch.  
*Advances in neural information processing systems*, 30, 2017.

- 648 Laixi Shi, Gen Li, Yuting Wei, Yuxin Chen, and Yuejie Chi. Pessimistic q-learning for offline  
649 reinforcement learning: Towards optimal sample complexity. In *International conference on*  
650 *machine learning*, pp. 19967–20025. PMLR, 2022.
- 651 I. Shpitser, T.J. VanderWeele, and J.M. Robins. On the validity of covariate adjustment for estimat-  
652 ing causal effects. In *Proceedings of the Twenty-Sixth Conference on Uncertainty in Artificial*  
653 *Intelligence*, pp. 527–536. AUAJ, Corvallis, OR, 2010.
- 654 Jörg Stoye. More on confidence intervals for partially identified parameters. *Econometrica*, 77(4):  
655 1299–1315, 2009.
- 656 Alexander L Strehl and Michael L Littman. An analysis of model-based interval estimation for  
657 markov decision processes. *Journal of Computer and System Sciences*, 74(8):1309–1331, 2008.
- 658 Richard S Sutton. Learning to predict by the methods of temporal differences. *Machine learning*, 3:  
659 9–44, 1988.
- 660 Richard S. Sutton and Andrew G. Barto. *Reinforcement learning: An introduction*. MIT press, 1998.
- 661 Adith Swaminathan and Thorsten Joachims. Counterfactual risk minimization: Learning from logged  
662 bandit feedback. In *International Conference on Machine Learning*, pp. 814–823, 2015.
- 663 Aviv Tamar, Shie Mannor, and Huan Xu. Scaling up robust mdps using function approximation. In  
664 *International conference on machine learning*, pp. 181–189. PMLR, 2014.
- 665 Philip S Thomas, Georgios Theodorou, and Mohammad Ghavamzadeh. High-confidence off-policy  
666 evaluation. In *AAAI*, pp. 3000–3006, 2015.
- 667 J. Tian and J. Pearl. Probabilities of causation: Bounds and identification. Technical report, 2000.
- 668 D Todem, J Fine, and L Peng. A global sensitivity test for evaluating statistical hypotheses with  
669 nonidentifiable models. *Biometrics*, 66(2):558–566, 2010.
- 670 Mark Towers, Ariel Kwiatkowski, Jordan Terry, John U Balis, Gianluca De Cola, Tristan Deleu,  
671 Manuel Goulão, Andreas Kallinteris, Markus Krimmel, Arjun KG, et al. Gymnasium: A standard  
672 interface for reinforcement learning environments. *arXiv preprint arXiv:2407.17032*, 2024.
- 673 Stijn Vansteelandt, Els Goetghebeur, Michael G Kenward, and Geert Molenberghs. Ignorance and  
674 uncertainty regions as inferential tools in a sensitivity analysis. *Statistica Sinica*, pp. 953–979,  
675 2006.
- 676 Yue Wang and Shaofeng Zou. Online robust reinforcement learning with model uncertainty. *Advances*  
677 *in Neural Information Processing Systems*, 34:7193–7206, 2021.
- 678 Christopher JCH Watkins and Peter Dayan. Q-learning. *Machine learning*, 8(3-4):279–292, 1992.
- 679 Christopher John Cornish Hellaby Watkins. *Learning from delayed rewards*. PhD thesis, University  
680 of Cambridge England, 1989.
- 681 Wolfram Wiesemann, Daniel Kuhn, and Berç Rustem. Robust markov decision processes. *Mathe-*  
682 *matics of Operations Research*, 38(1):153–183, 2013.
- 683 Huan Xu and Shie Mannor. Distributionally robust markov decision processes. *Advances in Neural*  
684 *Information Processing Systems*, 23, 2010.
- 685 Pengqian Yu and Huan Xu. Distributionally robust counterpart in markov decision processes. *IEEE*  
686 *Transactions on Automatic Control*, 61(9):2538–2543, 2015.
- 687 Huan Zhang, Hongge Chen, Chaowei Xiao, Bo Li, Mingyan Liu, Duane Boning, and Cho-Jui  
688 Hsieh. Robust deep reinforcement learning against adversarial perturbations on state observations.  
689 *Advances in Neural Information Processing Systems*, 33:21024–21037, 2020.
- 690 Junzhe Zhang and Elias Bareinboim. Near-optimal reinforcement learning in dynamic treatment  
691 regimes. In *Advances in Neural Information Processing Systems*, pp. 13401–13411, 2019.

702 Junzhe Zhang and Elias Bareinboim. Bounding causal effects on continuous outcomes. In *Proceedings*  
703 *of the 35nd AAAI Conference on Artificial Intelligence, 2021.*  
704  
705 Junzhe Zhang, Jin Tian, and Elias Bareinboim. Partial counterfactual identification from observational  
706 and experimental data. In *International Conference on Machine Learning*, pp. 26548–26558.  
707 PMLR, 2022.  
708  
709  
710  
711  
712  
713  
714  
715  
716  
717  
718  
719  
720  
721  
722  
723  
724  
725  
726  
727  
728  
729  
730  
731  
732  
733  
734  
735  
736  
737  
738  
739  
740  
741  
742  
743  
744  
745  
746  
747  
748  
749  
750  
751  
752  
753  
754  
755

## A RELATED WORK

Our work builds upon the literature on the partial identification of causal effects, sensitivity analysis, and robust reinforcement learning from offline data.

**Partial Identification and Sensitivity Analysis** Seminal work of Manski (1990) developed the first bounds on causal effects in non-identifiable settings using observational data in the single-stage treatment model with contextual information (i.e., a contextual bandit model). These bounds were then expanded to the instrumental variable setting (Balke & Pearl, 1997; Imbens & Angrist, 1994) partially identify counterfactual probabilities of causation (Tian & Pearl, 2000). More recently, (Zhang & Bareinboim, 2021) improved the bounds for applicability to continuous outcomes. (Zhang et al., 2022) established a general framework for estimating bounds on interventional and counterfactual effects. While Zhang et al. (2022) develop informative bounds using both observational and experimental data, they focus on general counterfactual queries by discretizing the exogenous latent space, formulating bounds as polynomial programs over this discretization and a Bayesian framework to approximately estimate bounds using MCMC.

Sensitivity analysis attempts to provide intervals on causal effects by assuming the level of confounding, for example, via models such as Marginal Sensitivity analysis, which considers deviations in the propensity score in relation to the estimated propensity (Rosenbaum, 2005; Richardson et al., 2014; Todem et al., 2010; Vansteelandt et al., 2006; Kallus & Zhou, 2018; Kallus et al., 2019; Namkoong et al., 2020; Jesson et al., 2022; Bruns-Smith & Zhou, 2023; Kausik et al., 2024). Other approaches explore additional parametric assumptions about the structural functions, including linearity (Cinelli et al., 2019) and Lipschitz continuity (Khan et al., 2023). Our work does not rely on additional functional constraints on the underlying system dynamics. Instead, we focus on the settings of standard discrete Markov Decision Processes (MDPs) with an infinite horizon. We develop robust off-policy evaluation algorithms to estimate closed-form bounds over the discounted cumulative rewards of candidate policies from offline observational data contaminated with unobserved confounding bias.

**Robust Reinforcement Learning** Unlike planning in a standard MDP, robust reinforcement learning does not assume the parametrization of the transition probability function in the underlying model to be precisely determined. Instead, it is contained in a set of model parameters which is called the uncertainty set (Iyengar, 2005; Nilim & El Ghaoui, 2005; Xu & Mannor, 2010; Wiesemann et al., 2013; Yu & Xu, 2015; Mannor et al., 2016; Petrik & Russel, 2019). The goal of the agent is to learn a robust policy that performs the best under the worst possible case in the uncertainty set. Similar problems have been studied under the rubrics of safe policy learning (Thomas et al., 2015; Ghavamzadeh et al., 2016) or pessimistic reinforcement learning (Shi et al., 2022).<sup>2</sup>

Robust RL algorithms with provable guarantees have been proposed in tabular settings or under the assumptions of linear functions (Lim et al., 2013; Tamar et al., 2014; Roy et al., 2017; Badrinath & Kalathil, 2021; Wang & Zou, 2021). Combined with the computational framework of deep learning, robust RL algorithms have been extended to complex, high-dimensional domains (Pinto et al., 2017; Zhang et al., 2020). More recently, (Panaganti et al., 2022) proposed Robust Fitted Q-Iteration (RFQI) to learn the best possible robust policy from offline data with theoretical guarantees on the performance of the learned policy. Our work differs from robust RL methods since it does not require a pre-specified uncertainty set of model parameters. Instead, we construct the ignorance region over the underlying system dynamics from the confounded observational data using partial causal identification. Based on the learned uncertainty set, we then derived closed-form bounds over the value functions of the target policy. *To the best of our knowledge, this is the first work that develops off-policy algorithms using eligibility traces to obtain evaluations of candidate policies from biased offline data, possibly contaminated with unmeasured confounding or no-overlap, with provable guarantees on the convergence of learned evaluations.*

<sup>2</sup>Indeed, the idea of planning over a convex set of model parameters have been explored in online reinforcement learning. (Strehl & Littman, 2008) utilized an extended dynamic programming to learn an optimistic policy over a confidence set of models to balance the trade-off between exploration and exploitation.

## B PROOFS

This section provides proof of the main theoretical results provided in the paper.

**Theorem 1** (Causal Bellman Equation). *For an MDP environment  $M$  with reward  $Y_t \in [a, b] \subseteq \mathbb{R}$ , for any policy  $\pi(x | s)$ , its state value function  $V_\pi(s) \in [V_\pi(s), \bar{V}_\pi(s)]$  for every state  $s \in \mathcal{S}$ , where bounds  $V_\pi, \bar{V}_\pi$  are solutions given by the following dynamic programs,<sup>3</sup>*

$$\langle V_\pi(s), \bar{V}_\pi(s) \rangle = \sum_x P(x | s) \left( \pi(x | s) \left( \tilde{\mathcal{R}}(s, x) + \gamma \sum_{s', x'} \tilde{\mathcal{T}}(s, x, s') \langle V_\pi(s'), \bar{V}_\pi(s') \rangle \right) \right) \quad (10)$$

$$+ \pi(\neg x | s) \left( \langle a, b \rangle + \gamma \left\langle \min_{s'} V_\pi(s'), \max_{s'} \bar{V}_\pi(s') \right\rangle \right) \quad (11)$$

*Proof.* Following the Bellman equation (Bellman, 1966), the state value function at state  $s \in \mathcal{S}$  is given by

$$V_\pi(s) = \sum_x \pi(x | s) \left( \mathcal{R}(s, x) + \gamma \sum_{s'} \mathcal{T}(s, x, s') V_\pi(s') \right) \quad (20)$$

Among the above quantities, the reward function  $\mathcal{R}$  is bounded from the observational distribution (Manski, 1990) as follows,

$$\tilde{\mathcal{R}}(s, x) P(x | s) + aP(\neg x | s) \leq \mathcal{R}(s, x) \leq \tilde{\mathcal{R}}(s, x) P(x | s) + bP(\neg x | s) \quad (21)$$

where  $\tilde{\mathcal{R}}$  is the nominal reward function computed from the observational distribution and is defined in Eq. (9). Replacing the reward function  $\mathcal{R}$  in Eq. (20) with the above lower bound gives

$$V_\pi(s) \geq \sum_x \pi(x | s) \left( \tilde{\mathcal{R}}(s, x) P(x | s) + aP(\neg x | s) + \gamma \sum_{s'} \mathcal{T}(s, x, s') V_\pi(s') \right) + \sum_x b\pi(x | s)P(\neg x | s) \quad (22)$$

Similarly, the transition distribution  $\mathcal{T}$  can be bounded from the observational distribution (Manski, 1990),

$$\tilde{\mathcal{T}}(s, x, s') P(x | s) \leq \mathcal{T}(s, x, s') \leq \tilde{\mathcal{T}}(s, x, s') P(x | s) + P(\neg x | s) \quad (23)$$

and  $\tilde{\mathcal{T}}$  is the nominal transition distribution computed from the observational distribution defined in Eq. (9). Minimizing the lower bound in Eq. (22) subject to the above observational constraints in Eq. (23) and  $\sum_{s'} \mathcal{T}(s, x, s') = 1$  gives the following lower bound:

$$V_\pi(s) \geq \sum_x \pi(x | s) P(x | s) \left( \tilde{\mathcal{R}}(s, x) + aP(\neg x | s) + \gamma \sum_{s'} \tilde{\mathcal{T}}(s, x, s') V_\pi(s') \right) + \sum_x \pi(x | s) P(\neg x | s) \left( b + \min_{s'} V_\pi(s') \right) \quad (24)$$

The above lower bound is achieved by setting the worst-case transition probability  $\mathcal{T}(s, x, s^*) = P(\neg x | s)$  for state  $s^* = \arg \min_{s'} V_\pi(s')$  and  $\mathcal{T}(s, x, s') = \tilde{\mathcal{T}}(s, x, s') P(x | s)$  for all the other state  $s' \neq s^*$ . Note that the second term of the above inequality could be further written as:

$$\sum_x \pi(x | s) P(\neg x | s) \left( a + \min_{s'} V_\pi(s') \right) \quad (25)$$

$$= \sum_x \pi(x | s) (1 - P(x | s)) \left( a + \min_{s'} V_\pi(s') \right) \quad (26)$$

$$= \sum_x \pi(x | s) \left( a + \min_{s'} V_\pi(s') \right) - \sum_x \pi(x | s) P(x | s) \left( a + \min_{s'} V_\pi(s') \right) \quad (27)$$

$$= \sum_x P(x | s) \left( a + \min_{s'} V_\pi(s') \right) - \sum_x \pi(x | s) P(x | s) \left( a + \min_{s'} V_\pi(s') \right) \quad (28)$$

<sup>3</sup> $\langle a, b \rangle$  is a vector containing a lower bound  $a$  and an upper bound  $b$ . We highlight quantities that are different from the standard Bellmen Equation.

The last step holds since for any constant real value  $C$ ,  $\sum_x \pi(x | s)C = \sum_x P(x | s)C$ . The above equation can be further written as

$$\sum_x \pi(x | s)P(\neg x | s) \left( a + \min_{s'} V_\pi(s') \right) = \sum_x \pi(\neg x | s)P(x | s) \left( a + \min_{s'} V_\pi(s') \right) \quad (29)$$

Replacing the second term in Eq. (24) gives

$$\begin{aligned} V_\pi(s) \geq \sum_x \pi(x | s)P(x | s) \left( \tilde{\mathcal{R}}(s, x) + bP(\neg x | s) + \gamma \sum_{s'} \tilde{\mathcal{T}}(s, x, s') V_\pi(s') \right) \\ + \sum_x \pi(\neg x | s)P(x | s) \left( a + \min_{s'} V_\pi(s') \right) \end{aligned} \quad (30)$$

After a few simplifications, we obtain

$$\begin{aligned} V_\pi(s) \geq P(x | s) \left( \pi(x | s) \left( \tilde{\mathcal{R}}(s, x) + \gamma \sum_{s', x'} \tilde{\mathcal{T}}(s, x, s') V_\pi(s') \right) \right. \\ \left. + \pi(\neg x | s) \left( a + \gamma \min_{s'} V_\pi(s') \right) \right) \end{aligned} \quad (31)$$

Finally, minimizing the value function  $V_\pi$  subject to the above inequality gives the lower bound  $\underline{V}_\pi$ . The upper bound  $\overline{V}_\pi$  over the state value function could be similarly derived.  $\square$

**Theorem 2** (Causal Bellman Equation). *For an MDP environment  $M$  with reward signals  $Y_t \in [a, b] \subseteq \mathbb{R}$ , for any policy  $\pi(x | s)$ , its state-action value function  $Q_\pi \in [Q_\pi(s, x), \overline{Q}_\pi(s, x)]$  for any state-action pair  $(s, x) \in \mathcal{S} \times \mathcal{X}$ , where bounds  $\underline{Q}_\pi, \overline{Q}_\pi$  are given by as follows,*

$$\langle Q_\pi(s, x), \overline{Q}_\pi(s, x) \rangle = P(x | s) \left( \tilde{\mathcal{R}}(s, x) + \gamma \sum_{s', x'} \tilde{\mathcal{T}}(s, x, s') \langle V_\pi(s'), \overline{V}_\pi(s') \rangle \right) \quad (12)$$

$$+ P(\neg x | s) \left( \langle a, b \rangle + \gamma \langle \min_{s'} \underline{V}_\pi(s'), \max_{s'} \overline{V}_\pi(s') \rangle \right) \quad (13)$$

*Proof.* Applying Bellman equation (Bellman, 1966) allows us to iteratively write the state-action value function for any state-action pair  $(s, x) \in \mathcal{S} \times \mathcal{X}$  as

$$Q_\pi(s, x) = \mathcal{R}(s, x) + \gamma \sum_{s'} \mathcal{T}(s, x, s') V_\pi(s') \quad (32)$$

where the reward function  $\mathcal{R}$  is bounded from the observational distribution (Manski, 1990) following Eq. (21). Replacing the reward function  $\mathcal{R}$  in the above equation with the corresponding lower bound gives

$$Q_\pi(s, x) \geq P(x | s) \left( \tilde{\mathcal{R}}(s, x) + \gamma \sum_{s'} \mathcal{T}(s, x, s') V_\pi(s') \right) + aP(\neg x | s) \quad (33)$$

Similarly, the transition distribution  $\mathcal{T}$  can be bounded from the observational distribution (Manski, 1990) following Eq. (23). Minimizing the lower bound in Eq. (33) subject to the above observational constraints in Eq. (23) and  $\sum_{s'} \mathcal{T}(s, x, s') = 1$  gives the following solution:

$$Q_\pi(s, x) \geq P(x | s) \left( \tilde{\mathcal{R}}(s, x) + \gamma \sum_{s'} \tilde{\mathcal{T}}(s, x, s') V_\pi(s') \right) + P(\neg x | s) \left( a + \min_{s'} V_\pi(s') \right) \quad (34)$$

This lower bound is achieved by setting the worst-case transition probability  $\mathcal{T}(s, x, s^*) = P(\neg x | s)$  for state  $s^* = \arg \min_{s'} V_\pi(s')$  and  $\mathcal{T}(s, x, s') = \tilde{\mathcal{T}}(s, x, s')P(x | s)$  for all the other state  $s' \neq s^*$ . Finally, notice that  $V_\pi(s)$  is a function of  $Q_\pi(s, x)$  and is given by  $V_\pi(s) = \sum_x \pi(x | s)Q_\pi(s, x)$ . Minimizing the state-action value function  $Q_\pi$  subject to the above inequality leads to the lower bound  $\underline{Q}_\pi$ . The upper bound  $\overline{Q}_\pi$  could be similarly derived.  $\square$

**Theorem 3.** For any behavior policy, for any choice of  $\lambda \in [0, 1]$  that does not depend on the actions chosen at each state, let parameters  $w$  and  $s^*$  be defined as follows: (1) Lower Bound  $\underline{V}_\pi$ :  $w = a$  and  $s^* = \arg \min_s V_t(s)$ ; (2) Upper Bound  $\overline{V}_\pi$ :  $w = b$  and  $s^* = \arg \max_s V_t(s)$ . Then, Alg. 1 with offline updating converges with probability 1 to lower bound  $\underline{V}_\pi$  and upper bound  $\overline{V}_\pi$ , respectively, under the usual step-size conditions on  $\alpha$ .

*Proof.* We will focus on the convergence of lower bound  $\underline{V}_\pi(s)$ ; the proof for the upper bound  $\overline{V}_\pi(s)$  follows analogously. The proof is structured in two stages. First, we consider the truncated lower bound estimates corresponding to Eq. (14), which sums the adjusted rewards obtained from the environment for only  $n$  steps, then uses the current estimate of the value function lower bound to approximate the remaining value:

$$\underline{R}_t^{(n)} = \sum_{k=0}^{n-1} \gamma^k \left( \pi_{t+k} y_{t+k} + \neg \pi_{t+k} (b + \gamma \min_{s'} V(s')) \right) \prod_{i=t}^{t+k-1} \pi_i + \gamma^n V(s_{t+n}) \prod_{i=t}^{t+k-1} \pi_i \quad (35)$$

We need to show that  $\underline{R}_t^{(n)} - \underline{V}_\pi$  is a contraction mapping in the max norm. If this is true for any  $n$ , then by applying the general convergence theorem, the  $n$ -step return converges to  $\underline{V}_\pi$ . Then any convex combination will also converge to  $\underline{V}_\pi$ . For example, any combination using a  $\lambda$  parameter in the style of eligibility traces will converge to  $\underline{V}_\pi$ .

The expected value of the adjusted return with regard to the observational distribution for state  $s$  can be expressed as follows <sup>4</sup>:

$$\mathbb{E} \left[ \underline{R}_t^{(n)} \mid S_t = s \right] \quad (36)$$

$$= \sum_{k=1}^n \sum_{\bar{s}_{1:k}, \bar{x}_{1:k}, \bar{y}_{1:k}} P(\bar{s}_{1:k}, \bar{x}_{1:k}, \bar{y}_{1:k}) \gamma^{k-1} \left( \pi_k y_k + \neg \pi_k (b + \min_{s'} V(s')) \right) \prod_{i=1}^{k-1} \pi_i \quad (37)$$

$$+ \sum_{\bar{s}_{1:n}, \bar{x}_{1:n}} P(\bar{s}_{1:n}, \bar{x}_{1:n}) \gamma^n V(s_n) \prod_{i=1}^{n-1} \pi_i \quad (38)$$

$$= \sum_{k=1}^n \gamma^{k-1} \sum_{\bar{s}_{1:k}, \bar{x}_{1:k}} \prod_{i=1}^{k-1} \tilde{T}(s_i, x_i, s_{i+1}) P(x_i \mid s_i) \pi(x_i \mid s_i) \quad (39)$$

$$\cdot \left( \pi(x_k \mid s_k) \tilde{\mathcal{R}}(s_k, x_k) + \neg \pi(x_k \mid s_k) (b + \gamma \min_{s'} V(s')) \right) \quad (40)$$

$$+ \gamma^n \sum_{\bar{s}_{1:n}, \bar{x}_{1:n}} \prod_{i=1}^{n-1} \tilde{T}(s_i, x_i, s_{i+1}) P(x_i \mid s_i) \pi(x_i \mid s_i) V(s_n) \quad (41)$$

By applying the extended Bellman equation for the lower bound  $\underline{V}_\pi$  iteratively  $n$  times, we obtain:

$$\underline{V}_\pi(s) = \sum_{k=1}^n \sum_{\bar{s}_{1:k}, \bar{x}_{1:k}} \gamma^{k-1} \prod_{i=1}^{k-1} \tilde{T}(s_i, x_i, s_{i+1}) P(x_i \mid s_i) \pi(x_i \mid s_i) \quad (42)$$

$$\cdot \left( \pi(x_k \mid s_k) \tilde{\mathcal{R}}(s_k, x_k) + \neg \pi(x_k \mid s_k) (b + \gamma \min_{s'} \underline{V}_\pi(s')) \right) \quad (43)$$

$$+ \gamma^n \sum_{\bar{s}_{1:n}, \bar{x}_{1:n}} \prod_{i=1}^{n-1} \tilde{T}(s_i, x_i, s_{i+1}) P(x_i \mid s_i) \pi(x_i \mid s_i) \underline{V}_\pi(s_n) \quad (44)$$

Therefore,

$$\max_s \left| \mathbb{E} \left[ \underline{R}_t^{(n)} \mid S_t = s \right] - \underline{V}_\pi(s) \right| \leq \gamma \max_s |V(s) - \underline{V}_\pi(s)| \quad (45)$$

This means that any  $n$ -step return is a contraction in the max norm, and therefore, by applying (Jaakkola et al., 1994, Theorem 1), it converges to  $\underline{V}_\pi(s)$ .

<sup>4</sup>We abuse notation a bit and ignore the expected value operator  $\mathbb{E}[\cdot]$  outside.

In the second stage, we show that by applying the updates of Alg. 1 for  $n$  successive steps, we perform the same update by using the  $n$ -step adjusted return  $\underline{R}_t^{(n)}$ . The eligibility trace for state  $s$  can be written as, for  $t_n \in \mathbf{t}(s)$ ,

$$e_t(s) = \gamma^{t-t_n} \prod_{i=t_n+1}^t \pi_i. \quad (46)$$

We have

$$\sum_{k=1}^n e_{t+k-1}(s) \delta_{t+k-1}(s) \quad (47)$$

$$= \sum_{k=1}^n \gamma^{k-1} \prod_{i=t+1}^{t+k-1} \pi_i \left( \pi_{t+k}(y_{t+k} + \gamma V(s_{t+k})) + \pi_{t+k} \left( b + \gamma \min_{s'} V(s') \right) \right) \quad (48)$$

$$- V(s_{t+k-1}) \quad (49)$$

$$= \sum_{k=0}^{n-1} \gamma^k \left( \pi_{t+k} y_{t+k} + \neg \pi_{t+k} \left( b + \gamma \min_{s'} V(s') \right) \right) \prod_{i=t}^{t+k-1} \pi_i + \gamma^n V(s_{t+n}) \prod_{i=t}^{t+k-1} \pi_i \quad (50)$$

$$- V(s_t) \quad (51)$$

$$= \underline{R}_t^{(n)} - V(s_t) \quad (52)$$

Since C-TD ( $\lambda$ ) is equivalent to applying a convex mixture of  $n$ -step updates, and each update converges to correct lower bounds  $\underline{V}_\pi$  for the state value functions, Alg. 1 converges to correct lower bounds as well.  $\square$

**Theorem 4.** For any behavior policy, for any choice of  $\lambda \in [0, 1]$  that does not depend on the actions chosen at each state, let parameters  $w$  and  $s^*$  be defined as follows: (1) Lower Bound  $\underline{Q}_\pi$ :  $w = a$  and  $s^* = \arg \min_s \sum_{x'} \pi(x' | s) Q_t(s, x')$ ; (2) Upper Bound  $\overline{Q}_\pi$ :  $w = b$  and  $s^* = \arg \max_s \sum_{x'} \pi(x' | s) Q_t(s, x')$ . Then, Alg. 2 with offline updating converges with probability 1 to lower bound  $\underline{Q}_\pi$  and upper bound  $\overline{Q}_\pi$ , respectively, under the usual step-size conditions on  $\alpha$ .

*Proof.* We will focus on the convergence of lower bound  $\underline{Q}_\pi(s, x)$ ; the proof for the upper bound  $\overline{Q}_\pi(s, x)$  follows analogously. This proof is structured in two stages. Let  $Q_n$  denote the  $n$ -step tree backup estimator defined in Eq. (19). First we show that  $\mathbb{E}[Q_n(s, x)] - \underline{Q}_\pi(s, x)$  is a contraction using a proof by induction.

Let  $Q$  be the current estimate of the lower bound for the value function. For  $n = 1$ ,

$$\max_{s,x} |\mathbb{E}[Q_1(s, x)] - \underline{Q}_\pi(s, x)| \quad (53)$$

$$= \max_{s,x} \left| P(x | s) \left( \tilde{\mathcal{R}}(s, x) + \gamma \sum_{s', x'} \tilde{\mathcal{T}}(s, x, s') \sum_{x'} \pi(x' | s') Q(s', x') \right) \right. \quad (54)$$

$$\left. + P(\neg x | s) \left( b + \gamma \min_{s'} \sum_{x'} \pi(x' | s') Q(s', x') \right) \right) \quad (55)$$

$$- P(x | s) \left( \tilde{\mathcal{R}}(s, x) + \gamma \sum_{s', x'} \tilde{\mathcal{T}}(s, x, s') \sum_{x'} \pi(x' | s') \underline{Q}_\pi(s', x') \right) \quad (56)$$

$$\left. - P(\neg x | s) \left( b + \gamma \min_{s'} \sum_{x'} \pi(x' | s') \underline{Q}_\pi(s', x') \right) \right| \quad (57)$$

$$\leq \gamma \max_{s,x} |Q(s, x) - \underline{Q}_\pi(s, x)| \quad (58)$$

For the induction step, we assume that

$$\max_{s,x} |\mathbb{E}[Q_n(s, x)] - \underline{Q}_\pi(s, x)| \leq \gamma \max_{s,x} |Q(s, x) - \underline{Q}_\pi(s, x)| \quad (59)$$

Next we want to show that the same holds for  $Q_{n+1}(s, x)$ . We can rewrite  $Q_{n+1}(s, x)$  as follows,

$$Q_{n+1}(s, x) = \mathbb{1}_{x_t=x} \left( y_t + \sum_{x'} \left( \mathbb{1}_{x' \neq x} \pi(x' | s_{t+1}) Q(s_{t+1}, x') + \mathbb{1}_{x'=x} Q_n(s_{t+1}, x) \right) \right) \quad (60)$$

$$+ \mathbb{1}_{x_t \neq x} \left( w + \sum_{x'} \pi(x' | s^*) Q(s^*, x') \right) \quad (61)$$

We must have

$$\max_{s,x} |\mathbb{E}[Q_{n+1}(s, x)] - \underline{Q}_\pi(s, x)| \quad (62)$$

$$= \max_{s,x} \left| P(x | s) \left( \tilde{\mathcal{R}}(s, x) + \gamma \sum_{s', x'} \tilde{\mathcal{T}}(s, x, s') \sum_{x'} \pi(x' | s') \right. \right. \quad (63)$$

$$\left. \mathbb{1}_{x' \neq x} Q(s', x') + \mathbb{1}_{x'=x} \mathbb{E}[Q_n(s', x)] \right) \quad (64)$$

$$+ P(\neg x | s) \left( b + \gamma \min_{s'} \sum_{x'} \pi(x' | s') Q(s', x') \right) \quad (65)$$

$$- P(x | s) \left( \tilde{\mathcal{R}}(s, x) + \gamma \sum_{s', x'} \tilde{\mathcal{T}}(s, x, s') \sum_{x'} \pi(x' | s') \underline{Q}_\pi(s', x') \right) \quad (66)$$

$$- P(\neg x | s) \left( b + \gamma \min_{s'} \sum_{x'} \pi(x' | s') \underline{Q}_\pi(s', x') \right) \quad (67)$$

$$\leq \gamma \max_{s,x} \left| P(x | s) \gamma \sum_{s', x'} \tilde{\mathcal{T}}(s, x, s') \sum_{x'} \pi(x' | s') \mathbb{1}_{x' \neq x} (Q(s', x') - \underline{Q}_\pi(s', x')) \right. \quad (68)$$

$$\left. + \mathbb{1}_{x'=x} \mathbb{E}[(Q_n(s', x) - \underline{Q}_\pi(s', x'))] \right) \quad (69)$$

$$+ P(\neg x | s) \min_{s'} \sum_{x'} \pi(x' | s') (Q(s', x') - \underline{Q}_\pi(s', x')) \quad (70)$$

$$\leq \gamma \max_{s,x} |Q(s, x) - \underline{Q}_\pi(s, x)| \quad (71)$$

By applying (Jaakkola et al., 1994, Theorem 1), we can conclude that any  $n$ -step adjusted return converges to the correct lower bound for the state-action value function. Since all the  $n$ -step returns converge to  $\underline{Q}_\pi$ , any convex linear combination of  $n$ -step returns also converges to  $\underline{Q}_\pi$ .

For the second part of the proof, we show that C-TB ( $\lambda$ ) with  $\lambda = 1$  for  $n$  steps is equivalent to using  $Q_n$ . The eligibility trace for a state-action pair  $(s, x)$  can be rewritten as:

$$e_t(s, x) = \gamma^k \prod_{i=t+1}^{t+k-1} \pi_{i+1} \mathbb{1}_{x_i=x}. \quad (72)$$

By adding and subtracting the weighted action value  $\pi_{t+k} \mathbb{1}_{x_{t+k}=x}$  for the action taken on each step from the return, and regrouping, we have

$$Q(s_t, x) + \sum_{k=1}^n \gamma^{k-1} \prod_{i=t+1}^{t+k-1} \pi_{i+1} \mathbb{1}_{x_i=x} \left( \mathbb{1}_{x_{t+k}=x} \left( y_{t+k} + \sum_{x' \neq x} \pi(x' | s_{t+k+1}) Q(s_{t+k+1}, x') \right) \right) \quad (73)$$

$$+ \mathbb{1}_{x_{t+k} \neq x} \left( w + \min_{s'} \sum_{x'} \pi(x' | s') Q(s', x') \right) - Q(s_{t+k}, x) \quad (74)$$

$$= Q(s_t, x) + \sum_{k=1}^n e_{t+k}(s_t, x) \delta_{t+k}(x) \quad (75)$$

This concludes the proof.  $\square$

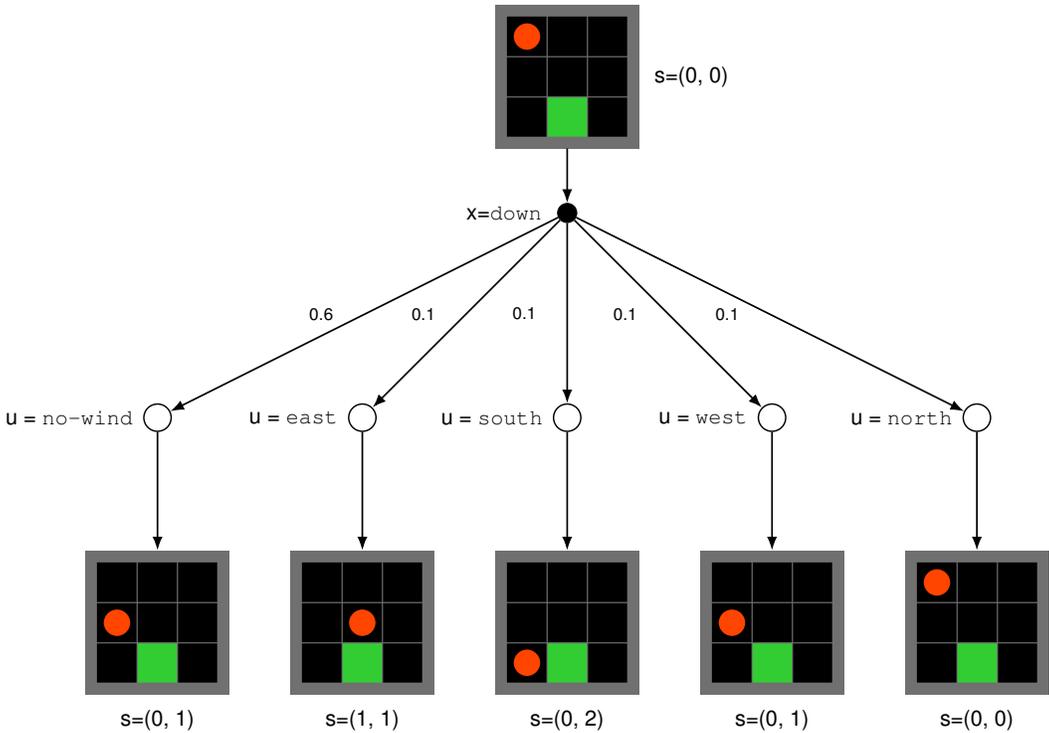


Figure 6: Trajectories sampled from the interventional transition distribution  $\mathcal{T}$ .

### C EXPERIMENTAL SETUPS

In this section, we provide details on the experimental setups and additional discussion on the simulation environment. All experiments were performed on a 2021 MacBook Pro with 16GB memory, implemented in Python. The simulation environment is built upon the Gymnasium framework (Brockman et al., 2016). We plan to release the source code with the camera-ready version of the manuscript.

**Windy Gridworld** Our simulation builds on the Windy Gridworld environment described in Fig. 1b, where the red dot represents the agent and the green square represents the goal state. The agent’s location is represented using a vector  $(i, j)$  where  $i \in \{0, 1, 2\}$  is the column index, and  $j \in \{0, 1, 2\}$  is the row index. So the agent’s starting state is  $(0, 0)$  and the goal state is  $(1, 2)$ . Fig. 7 shows the detailed state representation for each location in the gridworld.

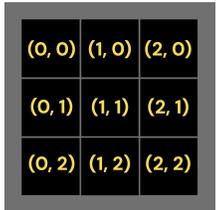
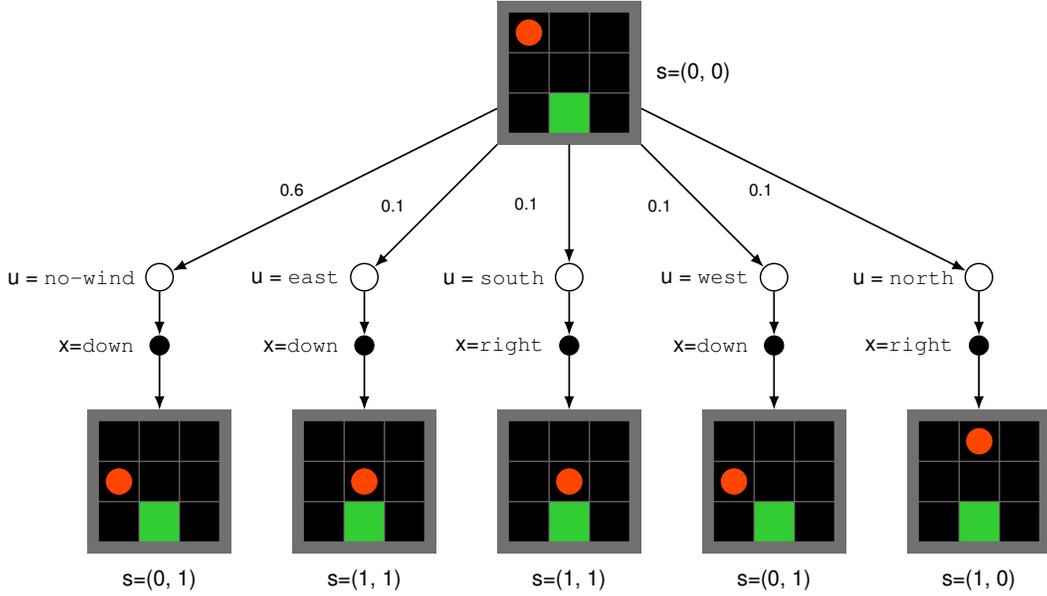


Figure 7: Agent’s state in Windy Gridworld environment.

The agent can take five actions  $x \in \mathcal{X}$  - up, down, right, left, and stay-put, corresponding to vector  $(0, -1), (0, 1), (1, 0), (-1, 0)$ , and  $(0, 0)$  respectively. Meanwhile, the agent’s movement is also affected by a wind; the wind direction  $u \in \mathcal{U}$  include - north, south, east, west, and no-wind, corresponding to vector  $(0, -1), (0, 1), (1, 0), (-1, 0)$ , and  $(0, 0)$  respectively. Table 1 summarizes the detailed parametrization for the agent’s action and the wind direction.

Action $x$	up	down	right	left	stay-put
Wind $u$	north	south	east	west	no-wind
Vector $v$	$(0, -1)$	$(0, 1)$	$(1, 0)$	$(-1, 0)$	$(0, 0)$

Table 1: Vector representations for the agent’s action  $X$  and the wind direction  $U$ .



1154 Figure 8: Trajectories sampled from the observational transition distribution  $\tilde{\mathcal{T}}$  induced by a con-  
1155 founded behavior policy  $f_X$ .

1156  
1157 every time step  $t = 1, 2, \dots$ , the wind  $U_t$  can blow in directions north, south, east, west with  
1158 equal probabilities of 10%; otherwise, the weather is nice and there is no-wind. That is,

$$1160 \forall i \in \{-1, 1\}, \quad P(U_t = (i, 0)) = P(U_t = (0, i)) = 0.1, \quad \text{and} \quad P(U_t = (0, 0)) = 0.6 \quad (76)$$

1161 At every time step  $t$ , the agent receives a constant reward  $Y_t \leftarrow -1$ . The next state of the agent is  
1162 shifted by both its action and the wind direction through the mechanism

$$1164 S_{t+1} \leftarrow \max \{ \min \{ S_t + X_t + U_t, (2, 2) \}, (0, 0) \}. \quad (77)$$

1165 In other words, the agent's next state  $S_{t+1}$  is a vector sum of the agent's current location  $S_t$ , its action  
1166  $X_t$ , and the wind direction  $U_t$ , truncated by the board's boundary  $i = 0, 2$  and  $j = 0, 2$ . For instance,  
1167 we show in Fig. 6 the system dynamics for the agent's interactions with the gridworld environment  
1168 at from the location  $s = (0, 0)$ , taking the action down ( $x = (0, 1)$ ). In this case, when the wind is  
1169 blowing towards south ( $u = (0, 1)$ ), the agent's location will be shifted by both the action  $x$  and the  
1170 windy direction  $u$ , and moves to the bottom left corner  $s' = (0, 2)$  at the next time step. Since among  
1171 all wind directions,  $u = \text{east}$  is the only latent state moving the agent to the center  $s' = (0, 2)$ , we  
1172 must have the following evaluation for the interventional distribution  $P_{X_t}(S_{t+1} | S_t)$ ,

$$1173 P_{X_t \leftarrow (0,1)}(S_{t+1} = (0, 2) | S_t = (0, 0)) = P(U_t = (1, 0)) \quad (78)$$

$$1174 = 0.1 \quad (79)$$

1175  
1176 That is, the agent's transition distribution  $\mathcal{T}(s, x, s') = 0.1$  when starting from  $s = (0, 1)$ , taking  
1177 action  $x = (0, 1)$ , and moving to the next state  $s' = (0, 2)$ .

1178  
1179 **Confounded Behavior Policy** Consider now an off-policy learning task in the windy gridworld,  
1180 where the agent's goal is to evaluate the expected return of a target policy  $\pi^*$  described in Fig. 2a.  
1181 Following such a policy  $\pi^*$ , the agent will consistently move towards the goal state  $s = (1, 2)$  from  
1182 its current location, regardless of the wind direction.

1183 The detailed parametrization of the agent's system dynamics in the windy gridworld remains unknown.  
1184 Instead, it has access to observed trajectories generated by a behavior policy  $x \leftarrow f_X(s, u)$  which  
1185 could sense the wind and select an action accordingly; Fig. 9 provides a detailed description for this  
1186 behavior policy. For example, when the agent is located in the top-left corner ( $s = (0, 0)$ ) and the  
1187 wind is blowing south ( $s = (0, 1)$ ), the behavior policy  $x \leftarrow f_X(s, u)$  will decide to move right  
( $x = (1, 0)$ ) so that the agent could get to the center ( $s' = (1, 1)$ ).

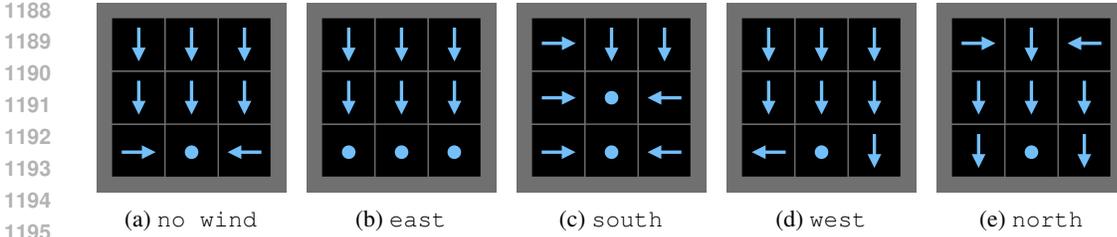


Figure 9: A confounded behavior policy  $f_X$  selecting values based on the agent’s location  $S$  and the latent wind direction  $U$ .

Consequently, the wind direction  $U_t$  becomes an unobserved confounder in the generative process for the offline observational data, affecting the allocated action  $X_t$  and the next state  $S_{t+1}$  simultaneously. The presence of unobserved confounders lead to violations of causal consistency (Def. 2). To witness, Fig. 8 shows observed trajectories in the offline data when the agent starts from state  $s = (0, 0)$ . When the weather is nice (no-wind) or the wind  $u$  is blowing towards east or west, the behavior policy selects action  $x = \text{down}$ , similar to the interventional trajectories of Fig. 6. On the other hand, when the wind is blowing towards north or south, the behavior policy selects action  $x = \text{right}$ , moving the agent towards the center of the board. Among all the possible next state in the observational data, we find that the agent will never reach the bottom left corner  $s = (0, 2)$ . This means that when evaluating the observational distribution  $P(S_{t+1} | S_t, X_t)$ , we must have

$$P(S_{t+1} = (0, 2) | S_t = (0, 0), X_t = (0, 1)) = 0 \quad (80)$$

In other words, the nominal transition distribution  $\tilde{T}(s, x, s') = 0$  when one observes the agent starting from  $s = (0, 1)$ , taking action  $x = (0, 1)$ , and moving to the next state  $s' = (0, 2)$ . Comparing the evaluations in Eqs. (79) and (80), we find that  $P_{x_t}(s_{t+1} | s_t) \neq P(s_{t+1} | s_t, x_t)$ , that is, causal consistency (Def. 2) does not hold between the agent’s system dynamics in windy gridworld and the observational distribution generated by the confounded behavior policy in Fig. 9.