
Provably Efficient Causal Reinforcement Learning with Confounded Observational Data

1 Empowered by neural networks, deep reinforcement learning (DRL) achieves tremendous empirical
2 successes. However, DRL requires a large dataset by interacting with the environment, which is un-
3 realistic in critical scenarios such as autonomous driving and personalized medicine. In this paper,
4 we study how to incorporate the dataset collected in the offline setting to improve the sample effi-
5 ciency in the online setting. To incorporate the observational data, we face two challenges. (a) The
6 behavior policy that generates the observational data may depend on unobserved random variables
7 (confounders), which affect the received rewards and transition dynamics. (b) Exploration in the
8 online setting requires quantifying the uncertainty given both the observational and interventional
9 data. To tackle such challenges, we propose the deconfounded optimistic value iteration (DOVI)
10 algorithm, which incorporates the confounded observational data in a provably efficient manner.
11 DOVI explicitly adjusts for the confounding bias in the observational data, where the confounders
12 are partially observed or unobserved. In both cases, such adjustments allow us to construct the bonus
13 based on a notion of information gain, which takes into account the amount of information acquired
14 from the offline setting. In particular, we prove that the regret of DOVI is smaller than the optimal
15 regret achievable in the pure online setting when the confounded observational data are informative
16 upon the adjustments.

17 1 Introduction

18 Empowered by the breakthrough in neural networks, deep reinforcement learning (DRL) achieves
19 significant empirical successes in various scenarios [19, 34, 23, 35]. Learning an expressive function
20 approximator necessitates collecting a large dataset. Specifically, in the online setting, it requires
21 the agent to interact with the environment for a large number of steps. For example, to learn a
22 human-level policy for playing Atari games, the agent has to interact with a simulator for more
23 than 10^8 steps [13]. However, in most scenarios, we do not have access to a simulator that allows
24 for trial and error without any cost. Meanwhile, in critical scenarios, e.g., autonomous driving and
25 personalized medicine, trial and error in the real world is unsafe and even unethical. As a result, it
26 remains challenging to apply DRL to more scenarios.

27 To bypass such a barrier, we study how to incorporate the dataset collected offline, namely the
28 observational data, to improve the sample efficiency of RL in the online setting [21]. In contrast
29 to the interventional data collected online in possibly expensive ways, observational data are often
30 abundantly available in various scenarios. For example, in autonomous driving, we have access
31 to trajectories generated by the drivers. As another example, in personalized medicine, we have
32 access to electronic health records from doctors. However, to incorporate the observational data in
33 a provably efficient way, we have to address two challenges.

- 34 • The observational data are possibly confounded. Specifically, there often exist unobserved random
35 variables, namely confounders, that causally affect the agent and the environment at the same
36 time. In particular, the policy used to generate the observational data, namely the behavior policy,
37 possibly depends on the confounders. Meanwhile, the confounders possibly affect the received
38 rewards and the transition dynamics.

39 In the example of autonomous driving [9, 22], the drivers may be affected by complicated traffic
40 or poor road design, resulting in traffic accidents even without misconduct. The complicated

41 traffic and poor road design subsequently affect both the action of the drivers and the outcome.
42 Therefore, it is unclear from the observational data whether the accidents are due to the actions
43 adopted by the drivers. Agents trained with such observational data may be unwilling to take any
44 actions under complicated traffic, jeopardizing the safety of passengers.

45 In the example of personalized medicine [28, 8], the patients may not be compliant with pre-
46 scriptions and instructions, which subsequently affects both the treatment and the outcome. As
47 another example, the doctor may prescribe medicine to patients based on patients’ socioeconomic
48 status (which could be inferred by the doctor through interacting with the patients). Meanwhile,
49 socioeconomic status affects the patients’ health condition and subsequently plays the role of the
50 confounder. In both scenarios, such confounders may be unavailable due to privacy or ethical con-
51 cerns. Such a confounding issue makes the observational data uninformative and even misleading
52 for identifying and estimating the causal effect, which is crucial for decision-making in the online
53 setting. In all the examples, it is unclear from the observational data whether the outcome is due
54 to the actions adopted.

55 • Even without the confounding issue, it remains unclear how the observational data may facilitate
56 exploration in the online setting, which is the key to the sample efficiency of RL. At the core of
57 exploration is uncertainty quantification. Specifically, quantifying the uncertainty that remains
58 given the dataset collected up to the current step, including the observational data and the inter-
59 ventional data, allows us to construct a bonus. When incorporated into the reward, such a bonus
60 encourages the agent to explore the less visited state-action pairs with more uncertainty. In par-
61 ticular, constructing such a bonus requires quantifying the amount of information carried over by
62 the observational data from the offline setting, which also plays a key role in characterizing the
63 regret, especially how much the observational data may facilitate reducing the regret.

64 Uncertainty quantification becomes even more challenging when the observational data are con-
65 founded. Specifically, as the behavior policy depends on the confounders, there is a mismatch
66 between the data generating processes in the offline setting and the online setting. As a result,
67 it remains challenging to quantify how much information carried over from the offline setting is
68 useful for the online setting, as the observational data are uninformative and even misleading due
69 to the confounding issue.

70 **Contribution.** To study causal reinforcement learning, we propose a class of Markov decision
71 processes (MDPs), namely confounded MDPs, which captures the data generating processes in both
72 the offline setting and the online setting as well as their mismatch due to the confounding issue.
73 In particular, we study two tractable cases of confounded MDPs in the episodic setting with linear
74 function approximation [40, 41, 16, 7].

75 • In the first case, the confounders are partially observed in the observational data. Assuming that
76 an observed subset of the confounders satisfies the backdoor criterion [30], we propose the decon-
77 founded optimistic value iteration (DOVI) algorithm, which explicitly corrects for the confound-
78 ing bias in the observational data using the backdoor adjustment.

79 • In the second case, the confounders are unobserved in the observational data. Assuming that there
80 exists an observed set of intermediate states that satisfies the frontdoor criterion [30], we propose
81 an extension of DOVI, namely DOVI⁺, which explicitly corrects for the confounding bias in the
82 observational data using the composition of two backdoor adjustments. We remark that DOVI⁺
83 follows the same principle of design as DOVI and defer the discussion of DOVI⁺ to §A.

84 In both cases, the adjustments allow DOVI and DOVI⁺ to incorporate the observational data into the
85 interventional data while bypassing the confounding issue. It further enables estimating the causal
86 effect of a policy on the received rewards and the transition dynamics with enlarged effective sample
87 size. Moreover, such adjustments allow us to construct the bonus based on a notion of information
88 gain, which takes into account the amount of information carried over from the offline setting.

89 In particular, we prove that DOVI and DOVI⁺ attain the $\Delta_H \cdot \sqrt{d^3 H^3 T}$ -regret up to logarithmic
90 factors, where d is the dimension of features, H is the length of each episode, and $T = HK$

91 is the number of steps taken in the online setting, where K is the number of episodes. Here the
 92 multiplicative factor $\Delta_H > 0$ depends on d , H , and a notion of information gain that quantifies the
 93 amount of information obtained from the interventional data additionally when given the properly
 94 adjusted observational data. When the observational data are unavailable or uninformative upon the
 95 adjustments, Δ_H is a logarithmic factor. Correspondingly, DOVI and DOVI⁺ attain the optimal
 96 \sqrt{T} -regret achievable in the pure online setting [40, 41, 16, 7]. When the observational data are
 97 sufficiently informative upon the adjustments, Δ_H decreases towards zero as the effective sample
 98 size of the observational data increases, which quantifies how much the observational data may
 99 facilitate exploration in the online setting.

100 **Related Work.** Our work is related to the study of causal bandit [20]. The goal of causal bandit is to
 101 obtain the optimal intervention in the online setting where the data generating process is described
 102 by a causal diagram. The previous study establishes causal bandit algorithms in the online setting
 103 [32, 25], the offline setting [17, 18], and a combination of both settings [11]. In contrast to this line
 104 of work, we study causal RL in a combination of the online setting and the offline setting. Causal
 105 RL is more challenging than causal bandit, which corresponds to $H = 1$, as it involves the transition
 106 dynamics and is more challenging in exploration. See §B for a detailed literature review on causal
 107 bandit.

108 Our work is related to the study of causal RL considered in various settings. [43] propose a model-
 109 based RL algorithm that solves dynamic treatment regimes (DTR), which involve a combination
 110 of the online setting and the offline setting. Their algorithm hinges on the analysis of sensitivity
 111 [26, 36, 4, 42], which constructs a set of feasible models of the transition dynamics based on the
 112 confounded observational data. Correspondingly, their algorithm achieves exploration by choosing
 113 an optimistic model of the transition dynamics from such a feasible set. In contrast, we propose a
 114 model-free RL algorithm, which achieves exploration through the bonus based on a notion of in-
 115 formation gain. It is worth mentioning that the assumption of [43] is weaker than ours as theirs
 116 does not allow for identifying the causal effect. As a result of partial identification, the regret of
 117 their algorithm is the same as the regret in the pure online setting as $T \rightarrow +\infty$. In contrast, our
 118 work instantiates the following framework in handling confounders for reinforcement learning. (a)
 119 First, we propose the estimation equation based on the observations, which identifies the causal ef-
 120 fect of actions on the cumulative reward. (b) Second, we conduct point estimation and uncertainty
 121 quantification based on observations and the estimation equation. (c) Finally, we conduct explo-
 122 ration based on the uncertainty quantification and achieve the regret reduction in the online setting.
 123 Consequently, the regret of our algorithm is smaller than the regret in the pure online setting by
 124 a multiplicative factor for all T . [24] propose a model-based RL algorithm in a combination of
 125 the online setting and the offline setting. Their algorithm uses a variational autoencoder (VAE) for
 126 estimating a structural causal model (SCM) based on the confounded observational data. In partic-
 127 ular, their algorithm utilizes the actor-critic algorithm to obtain the optimal policy in such an SCM.
 128 However, the regret of their algorithm remains unclear. [6] propose a model-based RL algorithm
 129 in the pure online setting that learns the optimal policy in a partially observable Markov decision
 130 process (POMDP). The regret of their algorithm also remains unclear. [33] utilize generative adver-
 131 sarial reinforcement learning to reconstruct transition dynamics with confounder, and [38] propose a
 132 model-based approach for POMDP based on adjustment with proxy variables. In contrast, our work
 133 utilizes backdoor and frontdoor adjustments to handle confounded observation.

134 2 Confounded Reinforcement Learning

135 **Structural Causal Model.** We denote a structural causal model (SCM) [30] by a tuple (A, B, F, P) .
 136 Here A is the set of exogenous (unobserved) variables, B is the set of endogenous (observed) vari-
 137 ables, F is the set of structural functions capturing the causal relations, which determines an en-
 138 dogenous variable $v \in B$ based on the other exogenous and endogenous variables, and P is the
 139 distribution of all the exogenous variables. We say that a pair of variables Y and Z are confounded
 140 by a variable W if they are both caused by W .

141 An intervention on a set of endogenous variables $X \subseteq B$ assigns a value x to X regardless of
 142 the other exogenous and endogenous variables as well as the structural functions. We denote by
 143 $\text{do}(X = x)$ the intervention on X and write $\text{do}(x)$ if it is clear from the context. Similarly, a
 144 stochastic intervention [27, 10] on a set of endogenous variables $X \subseteq B$ assigns a distribution p to
 145 X regardless of the other exogenous and endogenous variables as well as the structural functions.
 146 We denote by $\text{do}(X \sim p)$ the stochastic intervention on X .

147 **Confounded Markov Decision Process.** To characterize a Markov decision process (MDP) in the
 148 offline setting with observational data, which are possibly confounded, we introduce an SCM, where
 149 the endogenous variables are the states $\{s_h\}_{h \in [H]}$, actions $\{a_h\}_{h \in [H]}$, and rewards $\{r_h\}_{h \in [H]}$. Let
 150 $\{w_h\}_{h \in [H]}$ be the confounders. In §3, we assume that the confounders are partially observed, while
 151 in §A, we assume that they are unobserved. The set of structural functions F consists of the transition
 152 of states $s_{h+1} \sim \mathcal{P}_h(\cdot | s_h, a_h, w_h)$, the transition of confounders $w_h \sim \tilde{\mathcal{P}}_h(\cdot | s_h)$, the behavior
 153 policy $a_h \sim \nu_h(\cdot | s_h, w_h)$, which depends on the confounder w_h , and the reward function
 $r_h(s_h, a_h, w_h)$. See Figure 1 for the causal diagram that describes such an SCM.

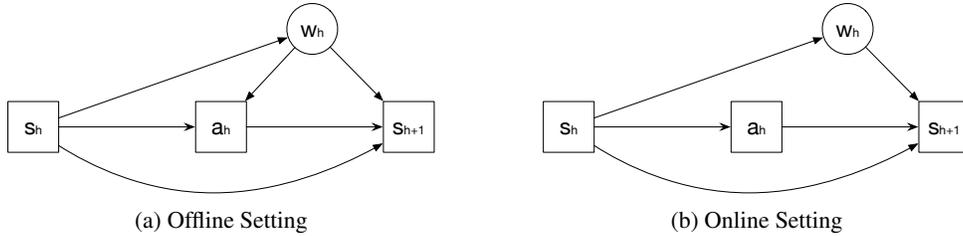


Figure 1: Causal diagrams of the h -th step of the confounded MDP (a) in the offline setting and (b) in the online setting, respectively.

154

155 Here a_h and s_{h+1} are confounded by w_h in addition to s_h . We denote such a confounded MDP
 156 by the tuple $(\mathcal{S}, \mathcal{A}, \mathcal{W}, H, \bar{\mathcal{P}}, r)$, where H is the length of an episode, \mathcal{S} , \mathcal{A} , and \mathcal{W} are the spaces
 157 of states, actions, and confounders, respectively, $r = \{r_h\}_{h \in [H]}$ is the set of reward functions,
 158 and $\bar{\mathcal{P}} = \{\mathcal{P}_h, \tilde{\mathcal{P}}_h\}_{h \in [H]}$ is the set of transition kernels. In the sequel, we assume without loss of
 159 generality that r_h takes value in $[0, 1]$ for all $h \in [H]$.

160 In the online setting that allows for intervention, we assume that the confounders $\{w_h\}_{h \in [H]}$
 161 are unobserved. A policy $\pi = \{\pi_h\}_{h \in [H]}$ induces the stochastic intervention $\text{do}(a_1 \sim$
 162 $\pi_1(\cdot | s_1), \dots, a_H \sim \pi_H(\cdot | s_H))$, which does not depend on the confounders. In particular, an
 163 agent interacts with the environment as follows. At the beginning of the k -th episode, the environment
 164 arbitrarily selects an initial state s_1^k and the agent selects a policy $\pi^k = \{\pi_h^k\}_{h \in [H]}$. At the
 165 h -th step of the k -th episode, the agent observes the state s_h^k and takes the action $a_h^k \sim \pi_h^k(\cdot | s_h^k)$.
 166 The environment randomly selects the confounder $w_h^k \sim \tilde{\mathcal{P}}_h(\cdot | s_h^k)$, which is unobserved, and the
 167 agent receives the reward $r_h^k = r_h(s_h^k, a_h^k, w_h^k)$. The environment then transits into the next state
 168 $s_{h+1}^k \sim \mathcal{P}_h(\cdot | s_h^k, a_h^k, w_h^k)$.

169 For a policy $\pi = \{\pi_h\}_{h \in [H]}$, which does not depend on the confounders $\{w_h\}_{h \in [H]}$, we define the
 170 value function $V^\pi = \{V_h^\pi\}_{h \in [H]}$ as follows,

$$V_h^\pi(s) = \mathbb{E}_\pi \left[\sum_{j=h}^H r_j(s_j, a_j, w_j) \mid s_h = s \right], \quad \forall h \in [H], \quad (2.1)$$

171 where we denote by \mathbb{E}_π the expectation with respect to the confounders $\{w_j\}_{j=h}^H$ and the trajectory
 172 $\{(s_j, a_j)\}_{j=h}^H$, starting from the state $s_j = s$ and following the policy π . Correspondingly, we define
 173 the action-value function $Q^\pi = \{Q_h^\pi\}_{h \in [H]}$ as follows,

$$Q_h^\pi(s, a) = \mathbb{E}_\pi \left[\sum_{j=h}^H r_j(s_j, a_j, w_j) \mid s_h = s, \text{do}(a_h = a) \right], \quad \forall h \in [H]. \quad (2.2)$$

174 We assess the performance of an algorithm using the regret against the globally optimal policy
 175 $\pi^* = \{\pi_h^*\}_{h \in [H]}$ in hindsight after K episodes, which is defined as follows,

$$\text{Regret}(T) = \max_{\pi} \sum_{k=1}^K (V_1^{\pi}(s_1^k) - V_1^{\pi^k}(s_1^k)) = \sum_{k=1}^K (V_1^{\pi^*}(s_1^k) - V_1^{\pi^k}(s_1^k)). \quad (2.3)$$

176 Here $T = HK$ is the total number of steps.

177 Our goal is to design an algorithm that minimizes the regret defined in (2.3), where π^* does not
 178 depend on the confounders $\{w_h\}_{h \in [H]}$. In the online setting that allows for intervention, it is well
 179 understood how to minimize such a regret [14, 3, 15, 16]. However, it remains unclear how to effi-
 180 ciently utilize the observational data obtained in the offline setting, which are possibly confounded.
 181 In real-world applications, e.g., autonomous driving and personalized medicine, such observational
 182 data are often abundant, whereas intervention in the online setting is often restricted. We refer to §D
 183 for a comparison between the confounded MDP and other extensions of MDP, including the dynam-
 184 ics treatment regime (DTR), partially observable MDP (POMDP), and contextual MDP (CMDP).

185 **Why is Incorporating Confounded Observational Data Challenging?** Straightforwardly incor-
 186 porating the confounded observational data into an online algorithm possibly leads to an undesirable
 187 regret due to the mismatch between the online and offline data generating processes. In particular,
 188 due to the existence of the confounders $\{w_h\}_{h \in [H]}$, which are partially observed (§3) or unobserved
 189 (§A), the conditional probability $\mathbb{P}(s_{h+1} | s_h, a_h)$ in the offline setting is different from the causal
 190 effect $\mathbb{P}(s_{h+1} | s_h, \text{do}(a_h))$ in the online setting [31]. More specifically, it holds that

$$\mathbb{P}(s_{h+1} | s_h, a_h) = \frac{\mathbb{E}_{w_h \sim \tilde{\mathcal{P}}_h(\cdot | s_h)} [\mathcal{P}_h(s_{h+1} | s_h, a_h, w_h) \cdot \nu_h(a_h | s_h, w_h)]}{\mathbb{E}_{w_h \sim \tilde{\mathcal{P}}_h(\cdot | s_h)} [\nu_h(a_h | s_h, w_h)]},$$

$$\mathbb{P}(s_{h+1} | s_h, \text{do}(a_h)) = \mathbb{E}_{w_h \sim \tilde{\mathcal{P}}_h(\cdot | s_h)} [\mathcal{P}_h(\cdot | s_h, a_h, w_h)].$$

191 In other words, without proper covariate adjustments [30], the confounded observational data may be
 192 not informative for estimating the transition dynamics and the associated action-value function in the
 193 online setting. To this end, we propose an algorithm that incorporates the confounded observational
 194 data in a provably efficient manner. Moreover, our analysis quantifies the amount of information
 195 carried over by the confounded observational data from the offline setting and to what extent it helps
 196 reducing the regret in the online setting.

197 3 Algorithm and Theory for Partially Observed Confounder

198 In this section, we propose the Deconfounded Optimistic Value Iteration (DOVI) algorithm. DOVI
 199 handles the case where the confounders are unobserved in the online setting but are partially ob-
 200 served in the offline setting. We then characterize the regret of DOVI. We defer the extension of
 201 DOVI, namely DOVI+, to §A which handles the case where the confounders are unobserved in both
 202 the online setting and the offline setting.

203 3.1 Algorithm

204 **Backdoor Adjustment.** In the online setting that allows for intervention, the causal effect of a_h on
 205 s_{h+1} given s_h , that is, $\mathbb{P}(s_{h+1} | s_h, \text{do}(a_h))$, plays a key role in the estimation of the action-value
 206 function. Meanwhile, the confounded observational data may not allow us to identify the causal
 207 effect $\mathbb{P}(s_{h+1} | s_h, \text{do}(a_h))$ if the confounder w_h is unobserved. However, if the confounder w_h is
 208 partially observed in the offline setting, the observed subset u_h of w_h allows us to identify the causal
 209 effect $\mathbb{P}(s_{h+1} | s_h, \text{do}(a_h))$, as long as u_h satisfies the following backdoor criterion.

210 **Assumption 3.1** (Backdoor Criterion [30, 31]). In the SCM defined in §2 and its induced directed
 211 acyclic graph (DAG), for all $h \in [H]$, there exists an observed subset u_h of w_h that satisfies the
 212 backdoor criterion, that is,

- 213 • the elements of u_h are not the descendants of a_h , and

214 • conditioning on s_h , the elements of u_h d -separate every path between a_h and s_{h+1} that has
 215 an incoming arrow into a_h .

216 See Figure 4 for an example that satisfies the backdoor criterion. In particular, we identify the causal
 217 effect $\mathbb{P}(s_{h+1} \mid s_h, \text{do}(a_h))$ as follows.

218 **Proposition 3.2** (Backdoor Adjustment [30]). Under Assumption 3.1, it holds for all $h \in [H]$ that

$$\begin{aligned} \mathbb{P}(s_{h+1} \mid s_h, \text{do}(a_h)) &= \mathbb{E}_{u_h \sim \mathbb{P}(\cdot \mid s_h)} [\mathbb{P}(s_{h+1} \mid s_h, a_h, u_h)], \\ \mathbb{E}[r_h(s_h, a_h, w_h) \mid s_h, \text{do}(a_h)] &= \mathbb{E}_{u_h \sim \mathbb{P}(\cdot \mid s_h)} [\mathbb{E}[r_h(s_h, a_h, w_h) \mid s_h, a_h, u_h]]. \end{aligned}$$

219 Here (s_{h+1}, s_h, a_h, u_h) follows the SCM defined in §2, which generates the confounded observa-
 220 tional data.

221 *Proof.* See [30] for a detailed proof. □

222 With a slight abuse of notation, we write $\mathbb{P}(s_{h+1} \mid s_h, a_h, u_h)$ as $\mathcal{P}_h(s_{h+1} \mid s_h, a_h, u_h)$ and
 223 $\mathbb{P}(u_h \mid s_h)$ as $\tilde{\mathcal{P}}_h(u_h \mid s_h)$, since they are induced by the SCM defined in §2. In the sequel, we
 224 define \mathcal{U} the space of observed state u_h and write $r_h = r_h(s_h, a_h, w_h)$ for notational simplicity.

225 **Backdoor-Adjusted Bellman Equation.** We now formulate the Bellman equation for the con-
 226 founded MDP. It holds for all $(s_h, a_h) \in \mathcal{S} \times \mathcal{A}$ that

$$Q_h^\pi(s_h, a_h) = \mathbb{E}_\pi \left[\sum_{j=h}^H r_j(s_j, a_j, u_j) \mid s_h, \text{do}(a_h) \right] = \mathbb{E}[r_h \mid s_h, \text{do}(a_h)] + \mathbb{E}_{s_{h+1}} [V_{h+1}^\pi(s_{h+1})],$$

227 where $\mathbb{E}_{s_{h+1}}$ denotes the expectation with respect to $s_{h+1} \sim \mathbb{P}(\cdot \mid s_h, \text{do}(a_h))$. Here
 228 $\mathbb{E}[r_h \mid s_h, \text{do}(a_h)]$ and $\mathbb{P}(\cdot \mid s_h, \text{do}(a_h))$ are characterized in Proposition 3.2. In the sequel, we define
 229 the following transition operator and counterfactual reward function,

$$(\mathbb{P}_h V)(s_h, a_h) = \mathbb{E}_{s_{h+1} \sim \mathbb{P}(\cdot \mid s_h, \text{do}(a_h))} [V(s_{h+1})], \quad \forall V : \mathcal{S} \mapsto \mathbb{R}, (s_h, a_h) \in \mathcal{S} \times \mathcal{A}, \quad (3.1)$$

$$R_h(s_h, a_h) = \mathbb{E}[r_h \mid s_h, \text{do}(a_h)], \quad \forall (s_h, a_h) \in \mathcal{S} \times \mathcal{A}. \quad (3.2)$$

230 We have the following Bellman equation,

$$Q_h^\pi(s_h, a_h) = R_h(s_h, a_h) + (\mathbb{P}_h V_{h+1}^\pi)(s_h, a_h), \quad \forall h \in [H], (s_h, a_h) \in \mathcal{S} \times \mathcal{A}. \quad (3.3)$$

231 Correspondingly, the Bellman optimality equation takes the following form,

$$Q_h^*(s_h, a_h) = R_h(s_h, a_h) + (\mathbb{P}_h V_{h+1}^*)(s_h, a_h), \quad V_h^*(s_h) = \max_{a_h \in \mathcal{A}} Q_h^*(s_h, a_h), \quad (3.4)$$

232 which holds for all $h \in [H]$ and $(s_h, a_h) \in \mathcal{S} \times \mathcal{A}$. Such a Bellman optimality equation allows us
 233 to adapt the least-squares value iteration (LSVI) algorithm [5, 14, 29, 3, 16].

234 **Linear Function Approximation.** We focus on the following setting with linear transition kernels
 235 and reward functions [40, 41, 16, 7], which corresponds to a linear SCM [31].

236 **Assumption 3.3** (Linear Confounded MDP). We assume that

$$\mathcal{P}_h(s_{h+1} \mid s_h, a_h, u_h) = \langle \phi_h(s_h, a_h, u_h), \mu_h(s_{h+1}) \rangle, \quad \forall h \in [H], (s_{h+1}, s_h, a_h) \in \mathcal{S} \times \mathcal{S} \times \mathcal{A},$$

237 where $\phi_h(\cdot, \cdot, \cdot)$ and $\mu_h(\cdot) = (\mu_{1,h}(\cdot), \dots, \mu_{d,h}(\cdot))^\top$ are \mathbb{R}^d -valued functions. We assume that
 238 $\sum_{i=1}^d \|\mu_{i,h}\|_1^2 \leq d$ and $\|\phi_h(s_h, a_h, u_h)\|_2 \leq 1$ for all $h \in [H]$ and $(s_h, a_h, u_h) \in \mathcal{S} \times \mathcal{A} \times \mathcal{U}$.
 239 Meanwhile, we assume that

$$\mathbb{E}[r_h \mid s_h, a_h, u_h] = \phi_h(s_h, a_h, u_h)^\top \theta_h, \quad \forall h \in [H], (s_h, a_h, u_h) \in \mathcal{S} \times \mathcal{A} \times \mathcal{U}, \quad (3.5)$$

240 where $\theta_h \in \mathbb{R}^d$ and $\|\theta_h\|_2 \leq \sqrt{d}$ for all $h \in [H]$.

241 Such a linear setting generalizes the tabular setting where \mathcal{S} , \mathcal{A} , and \mathcal{U} are finite.

242 **Proposition 3.4.** We define the backdoor-adjusted feature as follows,

$$\psi_h(s_h, a_h) = \mathbb{E}_{u_h \sim \tilde{\mathcal{P}}_h(\cdot | s_h)} [\phi_h(s_h, a_h, u_h)], \quad \forall h \in [H], (s_h, a_h) \in \mathcal{S} \times \mathcal{A}. \quad (3.6)$$

243 Under Assumption 3.1, it holds that

$$\mathbb{P}(s_{h+1} | s_h, \text{do}(a_h)) = \langle \psi_h(s_h, a_h), \mu_h(s_{h+1}) \rangle, \quad \forall h \in [H], (s_{h+1}, s_h, a_h) \in \mathcal{S} \times \mathcal{S} \times \mathcal{A}.$$

244 Moreover, the action-value functions Q_h^π and Q_h^* are linear in the backdoor-adjusted feature ψ_h for
245 all π .

246 *Proof.* See §F.1 for a detailed proof. \square

247 Such an observation allows us to estimate the action-value function based on the backdoor-adjusted
248 features $\{\psi_h\}_{h \in [H]}$ in the online setting. See §E for a detailed discussion. In the sequel, we assume
249 that either the density of $\{\tilde{\mathcal{P}}_h(\cdot | s_h)\}_{h \in [H]}$ is known or the backdoor-adjusted feature $\{\psi_h\}_{h \in [H]}$ is
250 known.

251 In the sequel, we introduce the DOVI algorithm (Algorithm 1). Each iteration of DOVI consists of
252 two components, namely point estimation, where we estimate Q_h^* based on the confounded observa-
253 tional data and the interventional data, and uncertainty quantification, where we construct the upper
254 confidence bound (UCB) of the point estimator.

Algorithm 1 Deconfounded Optimistic Value Iteration (DOVI) for Confounded MDP

Require: Observational data $\{(s_h^i, a_h^i, u_h^i, r_h^i)\}_{i \in [n], h \in [H]}$, tuning parameters $\lambda, \beta > 0$, backdoor-
adjusted feature $\{\psi_h\}_{h \in [H]}$, which is defined in (3.6).

- 1: **Initialization:** Set $\{Q_h^0, V_h^0\}_{h \in [H]}$ as zero functions and V_{H+1}^k as a zero function for $k \in [K]$.
 - 2: **for** $k = 1, \dots, K$ **do**
 - 3: **for** $h = H, \dots, 1$ **do**
 - 4: Set $\omega_h^k \leftarrow \operatorname{argmin}_{\omega \in \mathbb{R}^d} \sum_{\tau=1}^{k-1} (r_h^\tau + V_{h+1}^\tau(s_{h+1}^\tau) - \omega^\top \psi_h(s_h^\tau, a_h^\tau))^2 + \lambda \|\omega\|_2^2 + L_h^k(\omega)$,
where L_h^k is defined in (3.8).
 - 5: Set $Q_h^k(\cdot, \cdot) \leftarrow \min\{\psi_h(\cdot, \cdot)^\top \omega_h^k + \Gamma_h^k(\cdot, \cdot), H - h\}$, where Γ_h^k is defined in (3.12).
 - 6: Set $\pi_h^k(\cdot | s_h) \leftarrow \operatorname{argmax}_{a_h \in \mathcal{A}} Q_h^k(s_h, a_h)$ for all $s_h \in \mathcal{S}$.
 - 7: Set $V_h^k(\cdot) \leftarrow \langle \pi_h^k(\cdot | \cdot), Q_h^k(\cdot, \cdot) \rangle_{\mathcal{A}}$.
 - 8: **end for**
 - 9: Obtain s_1^k from the environment.
 - 10: **for** $h = 1, \dots, H$ **do**
 - 11: Take $a_h^k \sim \pi_h^k(\cdot | s_h^k)$. Obtain $r_h^k = r_h(s_h^k, a_h^k, u_h^k)$ and s_{h+1}^k .
 - 12: **end for**
 - 13: **end for**
-

255 **Point Estimation.** To solve the Bellman optimality equation in (3.4), we minimize the empirical
256 mean-squared Bellman error as follows at each step,

$$\omega_h^k \leftarrow \operatorname{argmin}_{\omega \in \mathbb{R}^d} \sum_{\tau=1}^{k-1} (r_h^\tau + V_{h+1}^\tau(s_{h+1}^\tau) - \omega^\top \psi_h(s_h^\tau, a_h^\tau))^2 + \lambda \|\omega\|_2^2 + L_h^k(\omega), \quad h = H, \dots, 1, \quad (3.7)$$

257 where we set $V_{H+1}^k = 0$ for all $k \in [K]$ and V_{h+1}^τ is defined in Line 7 of Algorithm 1 for all
258 $(\tau, h) \in [K] \times [H - 1]$. Here k is the index of episode, $\lambda > 0$ is a tuning parameter, and L_h^k is a
259 regularizer, which is constructed based on the confounded observational data. More specifically, we
260 define

$$L_h^k(\omega) = \sum_{i=1}^n (r_h^i + V_{h+1}^k(s_{h+1}^i) - \omega^\top \phi_h(s_h^i, a_h^i, u_h^i))^2, \quad \forall (k, h) \in [K] \times [H], \quad (3.8)$$

261 which corresponds to the least-squares loss for regressing $r_h^i + V_{h+1}^k(s_{h+1}^i)$ against $\phi_h(s_h^i, a_h^i, u_h^i)$
262 for all $i \in [n]$. Here $\{(s_h^i, a_h^i, u_h^i, r_h^i)\}_{(i, h) \in [n] \times [H]}$ are the confounded observational data, where

263 $u_h^i \sim \tilde{\mathcal{P}}_h(\cdot | s_h^i)$, $s_{h+1}^i \sim \mathcal{P}_h(\cdot | s_h^i, a_h^i, u_h^i)$, and $a_h^i \sim \nu_h(\cdot | s_h^i, w_h^i)$ with $\nu = \{\nu_h\}_{h \in [H]}$ being the
 264 behavior policy. Here recall that, with a slight abuse of notation, we write $\mathbb{P}(s_{h+1} | s_h, a_h, u_h)$ as
 265 $\mathcal{P}_h(s_{h+1} | s_h, a_h, u_h)$ and $\mathbb{P}(u_h | s_h)$ as $\tilde{\mathcal{P}}_h(u_h | s_h)$, since they are induced by the SCM defined in
 266 §2.

267 The update in (3.7) takes the following explicit form,

$$\omega_h^k \leftarrow (\Lambda_h^k)^{-1} \left(\sum_{\tau=1}^{k-1} \psi_h(s_h^\tau, a_h^\tau) \cdot (V_{h+1}^k(s_{h+1}^\tau) + r_h^\tau) + \sum_{i=1}^n \phi_h(s_h^i, a_h^i, u_h^i) \cdot (V_{h+1}^k(s_{h+1}^i) + r_h^i) \right), \quad (3.9)$$

268 where

$$\Lambda_h^k = \sum_{\tau=1}^{k-1} \psi_h(s_h^\tau, a_h^\tau) \psi_h(s_h^\tau, a_h^\tau)^\top + \sum_{i=1}^n \phi_h(s_h^i, a_h^i, u_h^i) \phi_h(s_h^i, a_h^i, u_h^i)^\top + \lambda I. \quad (3.10)$$

269 **Uncertainty Quantification.** We now construct the UCB $\Gamma_h^k(\cdot, \cdot)$ of the point estimator $\psi_h(\cdot, \cdot)^\top \omega_h^k$
 270 obtained from (3.9), which encourages the exploration of the less visited state-action pairs. To this
 271 end, we employ the following notion of information gain to motivate the UCB,

$$\Gamma_h^k(s_h^k, a_h^k) \propto H(\omega_h^k | \xi_{k-1}) - H(\omega_h^k | \xi_{k-1} \cup \{(s_h^k, a_h^k)\}), \quad (3.11)$$

272 where $H(\omega_h^k | \xi_{k-1})$ is the differential entropy of the random variable ω_h^k given the data ξ_{k-1} . In
 273 particular, $\xi_{k-1} = \{(s_h^\tau, a_h^\tau, r_h^\tau)\}_{(\tau, h) \in [k-1] \times [H]} \cup \{(s_h^i, a_h^i, u_h^i, r_h^i)\}_{(i, h) \in [n] \times [H]}$ consists of the
 274 confounded observational data and the interventional data up to the $(k-1)$ -th episode. However, it
 275 is challenging to characterize the distribution of ω_h^k . To this end, we consider a Bayesian counterpart
 276 of the confounded MDP, where the prior of ω_h^k is $N(0, I/\lambda)$ and the residual of the regression
 277 problem in (3.7) is $N(0, 1)$. In such a “parallel” confounded MDP, the posterior of ω_h^k follows
 278 $N(\mu_{k,h}, (\Lambda_h^k)^{-1})$, where Λ_h^k is defined in (3.10) and $\mu_{k,h}$ coincides with the right-hand side of
 279 (3.9). Moreover, it holds for all $(s_h^k, a_h^k) \in \mathcal{S} \times \mathcal{A}$ that

$$H(\omega_h^k | \xi_{k-1}) = 1/2 \cdot \log \det((2\pi e)^d \cdot (\Lambda_h^k)^{-1}),$$

$$H(\omega_h^k | \xi_{k-1} \cup \{(s_h^k, a_h^k)\}) = 1/2 \cdot \log \det\left((2\pi e)^d \cdot (\Lambda_h^k + \psi_h(s_h^k, a_h^k) \psi_h(s_h^k, a_h^k)^\top)^{-1}\right).$$

280 Correspondingly, we employ the following UCB, which instantiates (3.11), that is,

$$\Gamma_h^k(s_h^k, a_h^k) = \beta \cdot \left(\log \det(\Lambda_h^k + \psi_h(s_h^k, a_h^k) \psi_h(s_h^k, a_h^k)^\top) - \log \det(\Lambda_h^k) \right)^{1/2} \quad (3.12)$$

281 for all $(s_h^k, a_h^k) \in \mathcal{S} \times \mathcal{A}$. Here $\beta > 0$ is a tuning parameter. We highlight that, although the
 282 information gain in (3.11) relies on the “parallel” confounded MDP, the UCB in (3.12), which is used
 283 in Line 5 of Algorithm 1, does not rely on the Bayesian perspective. Also, our analysis establishes
 284 the frequentist regret.

285 **Regularization with Observational Data: A Bayesian Perspective.** In the “parallel” confounded
 286 MDP, it holds that

$$\omega_h^k \sim N(0, I/\lambda), \quad \omega_h^k | \xi_0 \sim N(\mu_{1,h}, (\Lambda_h^1)^{-1}), \quad \omega_h^k | \xi_{k-1} \sim N(\mu_{k,h}, (\Lambda_h^k)^{-1}),$$

287 where $\mu_{k,h}$ coincides with the right-hand side of (3.9) and $\mu_{1,h}$ is defined by setting $k = 1$ in
 288 $\mu_{k,h}$. Here $\xi_0 = \{(s_h^i, a_h^i, u_h^i, r_h^i)\}_{(i, h) \in [n] \times [H]}$ are the confounded observational data. Hence, the
 289 regularizer L_h^k in (3.8) corresponds to using $\omega_h^k | \xi_0$ as the prior for the Bayesian regression problem
 290 given only the interventional data $\xi_{k-1} \setminus \xi_0 = \{(s_h^\tau, a_h^\tau, r_h^\tau)\}_{(\tau, h) \in [k-1] \times [H]}$.

291 3.2 Theory

292 The following theorem characterizes the regret of DOVI, which is defined in (2.3).

293 **Theorem 3.5** (Regret of DOVI). Let $\beta = CdH\sqrt{\log(d(T+nH)/\zeta)}$ and $\lambda = 1$, where $C > 0$ and
 294 $\zeta \in (0, 1]$ are absolute constants. Under Assumptions 3.1 and 3.3, it holds with probability at least
 295 $1 - 5\zeta/2$ that

$$\text{Regret}(T) \leq C' \cdot \Delta_H \cdot \sqrt{d^3 H^3 T} \cdot \sqrt{\log(d(T+nH)/\zeta)}, \quad (3.13)$$

296 where $C' > 0$ is an absolute constant and

$$\Delta_H = \frac{1}{\sqrt{dH^2}} \sum_{h=1}^H (\log \det(\Lambda_h^{K+1}) - \log \det(\Lambda_h^1))^{1/2}. \quad (3.14)$$

297 *Proof.* See §F.3 for a detailed proof. \square

298 Note that $\Lambda_h^{K+1} \preceq (n+K+\lambda)I$ and $\Lambda_h^1 \succeq \lambda I$ for all $h \in [H]$. Hence, it holds that $\Delta_H =$
 299 $\mathcal{O}(\sqrt{\log(n+K+1)})$ in the worst case. Thus, the regret of DOVI is $\mathcal{O}(\sqrt{d^3 H^3 T})$ up to logarithmic
 300 factors, which is optimal in the total number of steps T if we only consider the online setting.
 301 However, Δ_H is possibly much smaller than $\mathcal{O}(\sqrt{\log(n+K+1)})$, depending on the amount of
 302 information carried over by the confounded observational data from the offline setting, which is
 303 quantified in the following.

304 **Interpretation of Δ_H : An Information-Theoretic Perspective.** Let ω_h^* be the parameter of the
 305 globally optimal action-value function Q_h^* , which corresponds to π^* in (2.3). Recall that we de-
 306 note by ξ_0 and ξ_K the confounded observational data $\{(s_h^i, a_h^i, u_h^i, r_h^i)\}_{(i,h) \in [n] \times [H]}$ and the union
 307 $\{(s_h^i, a_h^i, u_h^i, r_h^i)\}_{(i,h) \in [n] \times [H]} \cup \{(s_h^k, a_h^k, r_h^k)\}_{(k,h) \in [K] \times [H]}$ of the confounded observational data
 308 and the interventional data up to the K -th episode, respectively. We consider the aforementioned
 309 Bayesian counterpart of the confounded MDP, where the prior of ω_h^* is also $N(0, I/\lambda)$. In such a
 310 “parallel” confounded MDP, we have

$$\omega_h^* \sim N(0, I/\lambda), \quad \omega_h^* | \xi_0 \sim N(\mu_{1,h}^*, (\Lambda_h^1)^{-1}), \quad \omega_h^* | \xi_K \sim N(\mu_{K,h}^*, (\Lambda_h^{K+1})^{-1}), \quad (3.15)$$

311 where

$$\begin{aligned} \mu_{1,h}^* &= (\Lambda_h^1)^{-1} \sum_{i=1}^n \phi_h(s_h^i, a_h^i, u_h^i) \cdot (V_{h+1}^*(s_{h+1}^i) + r_h^i), \\ \mu_{K,h}^* &= (\Lambda_h^{K+1})^{-1} \left(\Lambda_h^1 \mu_{1,h}^* + \sum_{\tau=1}^K \psi_h(s_h^\tau, a_h^\tau) \cdot (V_{h+1}^*(s_{h+1}^\tau) + r_h^\tau) \right). \end{aligned}$$

312 It then holds for the right-hand side of (3.14) that

$$1/2 \cdot \log \det(\Lambda_h^{K+1}) - 1/2 \cdot \log \det(\Lambda_h^1) = H(\omega_h^* | \xi_0) - H(\omega_h^* | \xi_K). \quad (3.16)$$

313 The left-hand side of (3.16) characterizes the information gain of intervention in the online setting
 314 given the confounded observational data in the offline setting. In other words, if the confounded
 315 observational data are sufficiently informative upon the backdoor adjustment, then Δ_H is small,
 316 which implies that the regret is small. More specifically, the matrices $(\Lambda_h^1)^{-1}$ and $(\Lambda_h^{K+1})^{-1}$ de-
 317 fined in (3.10) characterize the ellipsoidal confidence sets given ξ_0 and ξ_K , respectively. If the
 318 confounded observational data are sufficiently informative upon the backdoor adjustment, Λ_h^{K+1}
 319 is close to Λ_h^1 . To illustrate, let $\{\psi_h(s_h^\tau, a_h^\tau)\}_{(\tau,h) \in [K] \times [H]}$ and $\{\phi_h(s_h^i, a_h^i, u_h^i)\}_{(i,h) \in [n] \times [H]}$
 320 be sampled uniformly at random from the canonical basis $\{e_\ell\}_{\ell \in [d]}$ of \mathbb{R}^d . It then holds that
 321 $\Lambda_h^{K+1} \approx (K+n)I/d + \lambda I$ and $\Lambda_h^1 \approx nI/d + \lambda I$. Hence, for $\lambda = 1$ and sufficiently large n and
 322 K , we have $\Delta_H = \mathcal{O}(\sqrt{\log(1+K/(n+d))}) = \mathcal{O}(\sqrt{K/(n+d)})$. For example, for $n = \Omega(K^2)$,
 323 it holds that $\Delta_H = \mathcal{O}(n^{-1/2})$, which implies that the regret of DOVI is $\mathcal{O}(n^{-1/2} \cdot \sqrt{d^3 H^3 T})$. In
 324 other words, if the confounded observational data are sufficiently informative upon the backdoor
 325 adjustment, the regret of DOVI can be arbitrarily small given a sufficiently large sample size n of
 326 the confounded observational data, which is often the case in practice [28, 8, 9, 22, 21].

327 **References**

- 328 [1] Abbasi-Yadkori, Y., Pál, D. and Szepesvári, C. (2011). Improved algorithms for linear stochastic
329 bandits. In *Advances in Neural Information Processing Systems*.
- 330 [2] Auer, P. and Ortner, R. (2007). Logarithmic online regret bounds for undiscounted reinforcement
331 learning. In *Advances in Neural Information Processing Systems*.
- 332 [3] Azar, M. G., Osband, I. and Munos, R. (2017). Minimax regret bounds for reinforcement
333 learning. In *International Conference on Machine Learning*.
- 334 [4] Balke, A. and Pearl, J. (2013). Counterfactuals and policy analysis in structural models. *arXiv
335 preprint arXiv:1302.4929*.
- 336 [5] Bradtke, S. J. and Barto, A. G. (1996). Linear least-squares algorithms for temporal difference
337 learning. *Machine Learning*, **22** 33–57.
- 338 [6] Buesing, L., Weber, T., Zwols, Y., Racaniere, S., Guez, A., Lespiau, J.-B. and Heess, N.
339 (2018). Woulda, coulda, shoulda: Counterfactually-guided policy search. *arXiv preprint
340 arXiv:1811.06272*.
- 341 [7] Cai, Q., Yang, Z., Jin, C. and Wang, Z. (2019). Provably efficient exploration in policy opti-
342 mization. *arXiv preprint arXiv:1912.05830*.
- 343 [8] Chakraborty, B. and Murphy, S. A. (2014). Dynamic treatment regimes. *Annual Review of
344 Statistics and Its Application*, **1** 447–464.
- 345 [9] de Haan, P., Jayaraman, D. and Levine, S. (2019). Causal confusion in imitation learning. In
346 *Advances in Neural Information Processing Systems*.
- 347 [10] Díaz, I. and Hejazi, N. (2019). Causal mediation analysis for stochastic interventions. *arXiv
348 preprint arXiv:1901.02776*.
- 349 [11] Forney, A., Pearl, J. and Bareinboim, E. (2017). Counterfactual data-fusion for online rein-
350 forcement learners. In *International Conference on Machine Learning*.
- 351 [12] Hallak, A., Di Castro, D. and Mannor, S. (2015). Contextual Markov decision processes. *arXiv
352 preprint arXiv:1502.02259*.
- 353 [13] Hessel, M., Modayil, J., Van Hasselt, H., Schaul, T., Ostrovski, G., Dabney, W., Horgan, D.,
354 Piot, B., Azar, M. and Silver, D. (2018). Rainbow: Combining improvements in deep rein-
355 forcement learning. In *AAAI Conference on Artificial Intelligence*.
- 356 [14] Jaksch, T., Ortner, R. and Auer, P. (2010). Near-optimal regret bounds for reinforcement learn-
357 ing. *Journal of Machine Learning Research*, **11** 1563–1600.
- 358 [15] Jin, C., Allen-Zhu, Z., Bubeck, S. and Jordan, M. I. (2018). Is Q-learning provably efficient?
359 In *Advances in Neural Information Processing Systems*.
- 360 [16] Jin, C., Yang, Z., Wang, Z. and Jordan, M. I. (2019). Provably efficient reinforcement learning
361 with linear function approximation. *arXiv preprint arXiv:1907.05388*.
- 362 [17] Kallus, N. and Zhou, A. (2018). Confounding-robust policy improvement. In *Advances in
363 Neural Information Processing Systems*.
- 364 [18] Kallus, N. and Zhou, A. (2018). Policy evaluation and optimization with continuous treat-
365 ments. *arXiv preprint arXiv:1802.06037*.
- 366 [19] Kober, J., Bagnell, J. A. and Peters, J. (2013). Reinforcement learning in robotics: A survey.
367 *International Journal of Robotics Research*, **32** 1238–1274.

- 368 [20] Lattimore, F., Lattimore, T. and Reid, M. D. (2016). Causal bandits: Learning good interven-
369 tions via causal inference. In *Advances in Neural Information Processing Systems*.
- 370 [21] Levine, S., Kumar, A., Tucker, G. and Fu, J. (2020). Offline reinforcement learning: Tutorial,
371 review, and perspectives on open problems. *arXiv preprint arXiv:2005.01643*.
- 372 [22] Li, C., Chan, S. H. and Chen, Y.-T. (2020). Who make drivers stop? Towards driver-
373 centric risk assessment: Risk object identification via causal inference. *arXiv preprint*
374 *arXiv:2003.02425*.
- 375 [23] Li, J., Monroe, W., Ritter, A., Galley, M., Gao, J. and Jurafsky, D. (2016). Deep reinforcement
376 learning for dialogue generation. *arXiv preprint arXiv:1606.01541*.
- 377 [24] Lu, C., Schölkopf, B. and Hernández-Lobato, J. M. (2018). Deconfounding reinforcement
378 learning in observational settings. *arXiv preprint arXiv:1812.10576*.
- 379 [25] Lu, Y., Meisami, A., Tewari, A. and Yan, Z. (2019). Regret analysis of causal bandit problems.
380 *arXiv preprint arXiv:1910.04938*.
- 381 [26] Manski, C. F. (1990). Nonparametric bounds on treatment effects. *American Economic Re-*
382 *view*, **80** 319–323.
- 383 [27] Muñoz, I. D. and van der Laan, M. (2012). Population intervention causal effects based on
384 stochastic interventions. *Biometrics*, **68** 541–549.
- 385 [28] Murphy, S. A. (2003). Optimal dynamic treatment regimes. *Journal of the Royal Statistical*
386 *Society: Series B (Statistical Methodology)*, **65** 331–355.
- 387 [29] Osband, I., Van Roy, B. and Wen, Z. (2014). Generalization and exploration via randomized
388 value functions. *arXiv preprint arXiv:1402.0635*.
- 389 [30] Pearl, J. (2009). *Causality*. Cambridge university press.
- 390 [31] Peters, J., Janzing, D. and Schölkopf, B. (2017). *Elements of Causal Inference: Foundations*
391 *and Learning Algorithms*. MIT press.
- 392 [32] Sen, R., Shanmugam, K., Dimakis, A. G. and Shakkottai, S. (2017). Identifying best interven-
393 tions through online importance sampling. In *International Conference on Machine Learning*.
- 394 [33] Shang, W., Yu, Y., Li, Q., Qin, Z., Meng, Y. and Ye, J. (2019). Environment reconstruction
395 with hidden confounders for reinforcement learning based recommendation. In *Proceedings*
396 *of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*.
- 397 [34] Silver, D., Huang, A., Maddison, C. J., Guez, A., Sifre, L., Van Den Driessche, G.,
398 Schrittwieser, J., Antonoglou, I., Panneershelvam, V., Lanctot, M. et al. (2016). Master-
399 ing the game of Go with deep neural networks and tree search. *Nature*, **529** 484.
- 400 [35] Silver, D., Schrittwieser, J., Simonyan, K., Antonoglou, I., Huang, A., Guez, A., Hubert, T.,
401 Baker, L., Lai, M., Bolton, A. et al. (2017). Mastering the game of Go without human knowl-
402 edge. *Nature*, **550** 354.
- 403 [36] Tan, Z. (2006). A distributional approach for causal inference using propensity scores. *Journal*
404 *of the American Statistical Association*, **101** 1619–1637.
- 405 [37] Tennenholtz, G., Mannor, S. and Shalit, U. (2019). Off-policy evaluation in partially observ-
406 able environments. *arXiv preprint arXiv:1909.03739*.
- 407 [38] Tennenholtz, G., Shalit, U. and Mannor, S. (2020). Off-policy evaluation in partially observ-
408 able environments. In *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 34.

- 409 [39] Vershynin, R. (2010). Introduction to the non-asymptotic analysis of random matrices. *arXiv preprint arXiv:1011.3027*.
410
- 411 [40] Yang, L. and Wang, M. (2019). Sample-optimal parametric Q-learning using linearly additive
412 features. In *International Conference on Machine Learning*.
- 413 [41] Yang, L. F. and Wang, M. (2019). Reinforcement learning in feature space: Matrix bandit,
414 kernels, and regret bound. *arXiv preprint arXiv:1905.10389*.
- 415 [42] Zhang, J. and Bareinboim, E. (2017). Transfer learning in multi-armed bandit: A causal ap-
416 proach. In *Autonomous Agents and Multi-Agent Systems*.
- 417 [43] Zhang, J. and Bareinboim, E. (2019). Near-optimal reinforcement learning in dynamic treat-
418 ment regimes. In *Advances in Neural Information Processing Systems*.

419 Checklist

- 420 1. For all authors...
- 421 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s
422 contributions and scope? [Yes]
- 423 (b) Did you describe the limitations of your work? [No]
- 424 (c) Did you discuss any potential negative societal impacts of your work? [No]
- 425 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
426 them? [Yes]
- 427 2. If you are including theoretical results...
- 428 (a) Did you state the full set of assumptions of all theoretical results? [Yes]
- 429 (b) Did you include complete proofs of all theoretical results? [Yes]
- 430 3. If you ran experiments...
- 431 (a) Did you include the code, data, and instructions needed to reproduce the main experi-
432 mental results (either in the supplemental material or as a URL)? [N/A]
- 433 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
434 were chosen)? [N/A]
- 435 (c) Did you report error bars (e.g., with respect to the random seed after running experi-
436 ments multiple times)? [N/A]
- 437 (d) Did you include the total amount of compute and the type of resources used (e.g., type
438 of GPUs, internal cluster, or cloud provider)? [N/A]
- 439 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 440 (a) If your work uses existing assets, did you cite the creators? [N/A]
- 441 (b) Did you mention the license of the assets? [N/A]
- 442 (c) Did you include any new assets either in the supplemental material or as a URL? [N/A]
- 443 (d) Did you discuss whether and how consent was obtained from people whose data
444 you’re using/curating? [N/A]
- 445 (e) Did you discuss whether the data you are using/curating contains personally identifi-
446 able information or offensive content? [N/A]
- 447 5. If you used crowdsourcing or conducted research with human subjects...
- 448 (a) Did you include the full text of instructions given to participants and screenshots, if
449 applicable? [N/A]
- 450 (b) Did you describe any potential participant risks, with links to Institutional Review
451 Board (IRB) approvals, if applicable? [N/A]
- 452 (c) Did you include the estimated hourly wage paid to participants and the total amount
453 spent on participant compensation? [N/A]