
Efficient RL with Impaired Observability: Learning to Act with Delayed and Missing State Observations

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 In real-world reinforcement learning (RL) systems, various forms of *impaired*
2 *observability* can complicate matters. These situations arise when an agent is
3 unable to observe the most recent state of the system due to latency or lossy
4 channels, yet the agent must still make real-time decisions. This paper introduces
5 a theoretical investigation into efficient RL in control systems where agents must
6 act with delayed and missing state observations. We establish near-optimal regret
7 bounds, of the form $\tilde{O}(\sqrt{\text{poly}(H)SAK})$, for RL in both the delayed and missing
8 observation settings. Despite impaired observability posing significant challenges
9 to the policy class and planning, our results demonstrate that learning remains
10 efficient, with the regret bound optimally depending on the state-action size of the
11 original system. Additionally, we provide a characterization of the performance of
12 the optimal policy under impaired observability, comparing it to the optimal value
13 obtained with full observability.

14 1 Introduction

15 In Reinforcement Learning (RL), an agent engages with an environment in a sequential manner. In
16 an ideal setting, at each time step, the agent would observe the current state of the environment, select
17 an action to perform, and receive a reward [Smallwood and Sondik, 1973, Bertsekas, 2012, Sutton
18 and Barto, 2018, Lattimore and Szepesvári, 2020]. However, real-world engineering systems often
19 introduce impaired observability and latency, where the agent may not have immediate access to the
20 instant state and reward information. In systems with lossy communication channels, certain state
21 observations may even be permanently missing, never reaching the agent. Nevertheless, the agent is
22 still required to make real-time decisions based on the available information.

23 The presence of impaired observability transforms the system into a complex interactive decision
24 process (Figure 1), presenting challenges for both learning and planning in RL. With limited knowl-
25 edge about recent states and rewards, the agent’s policy must extract information from the observed
26 history and utilize it to make immediate decisions. This introduces significant complexity to the
27 policy class and poses difficulties for RL. Moreover, the loss of information due to permanently
28 missing observations further hampers the efficiency of RL methods. Although a naïve approach
29 would involve augmenting the state and action space to create a fully observable Markov Decision
30 Process (MDP), such a method would lead to exponential regret growth in the state-action size.

31 **Why existing methods do not work.** One may be tempted to cast the problem of impaired
32 observability into a Partially Observed MDPs (POMDPs). However, this would not solve the
33 problem. In POMDP, the system does not reveal its instant state to the agent but provides an
34 emission state observation conditioned on the latent state. POMDPs are known to suffer from the
35 curse of history [Papadimitriou and Tsitsiklis, 1987, Bertsekas, 2012, Krishnamurthy, 2016], unless
36 additional assumptions are imposed. Existing efficient algorithms focus on subclasses of POMDPs

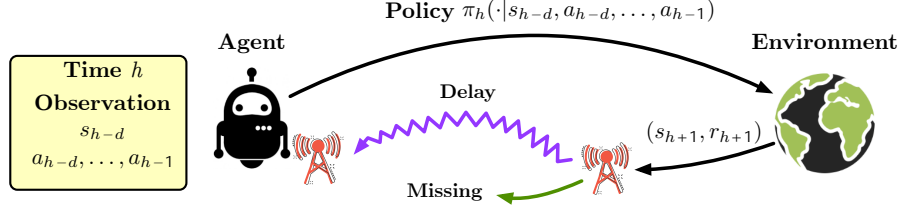


Figure 1: Reinforcement learning with impaired observability. At time h , the agent only observes the past state s_{h-d} and actions a_{h-d}, \dots, a_{h-1} . The policy depends on the observed information.

37 with decodable or distinguishable partial observations [Jin et al., 2020, Uehara et al., 2022, Zhan
 38 et al., 2022, Chen et al., 2022, Liu et al., 2022, Zhong et al., 2022, Chen et al., 2023], where the
 39 unseen instant state can be inferred from recent observations. Unfortunately, MDPs with impaired
 40 observability do not fall into these benign subclasses. The reason behind this is that at each time step,
 41 a new observation, if any, is in fact a past state. Viewing it as an emission state of the current one
 42 leads to a time reversal posterior distribution depending on the underlying transitions, which suffers
 43 from the curse of history and makes the POMDP intractable. The problem becomes even harder if
 44 some observations get missing.

45 Empirical evidences suggested that efficient RL is possible even with impaired state observability
 46 [Lizotte et al., 2008, Liu et al., 2014, Agarwal and Aggarwal, 2021]. However, theoretical under-
 47 standing of this problem is very limited. One notable work [Walsh et al., 2007] studied learning with
 48 constant-time delayed observations. They identified subclasses of MDPs with nearly deterministic
 49 transitions that can be efficiently learned. Beyond this special case, efficient RL with impaired
 50 observability in MDPs with fully generality remains largely open.

51 Some recent works studied delayed feedback in MDPs [Yang et al., 2023, Howson et al., 2023]. It
 52 is a fundamentally different problem where the agent’s policy can still access real-time states but
 53 learning uses delayed data. Our problem is fundamentally harder because the agent’s policy can only
 54 access the lossy and delayed history. See Section 1.1 for more discussions.

55 **Our results.** In this paper, we provide algorithms and regret analysis for learning the optimal policy
 56 in tabular MDPs with impaired observability. Note that this optimal policy is a different one from the
 57 optimal policy with full observability. To approach this problem, we construct an augmented MDP
 58 reformulation where the original state space is expanded to include available observations of past
 59 state and an action sequence. However, the expanded state space is much larger than the original one
 60 and naïve application of known methods would lead to exponentially large regret bounds. In our
 61 analysis, we exploit structure of the augmented transition model to achieve efficient learning and
 62 sharp regret bounds. The main results are summarized as follows.

63 • For MDPs with stochastic delays, we prove a sharp $\tilde{O}(H^4 \sqrt{SAK})$ regret bound (Theorem 4.1)
 64 comparing to the best feasible policy, Here S and A are the sizes of the original state and action
 65 spaces, respectively, H is the horizon, and K is the number of episodes. Here we allows the delay to
 66 be stochastic and conditionally independent given on current state and action. Moreover, we quantify
 67 the performance degradation of optimal value due to impaired observability, compared to optimal
 68 value of fully observable MDPs (Proposition B.2). We also showcase in Proposition 4.2 that a short
 69 delay does not reduce the optimal value, but slightly longer delay leads to substantial degradation.

70 • For MDPs with randomly missing observations, we provide an optimistic RL method that provably
 71 achieves $\tilde{O}(\sqrt{H^3 S^2 AK})$ regret (Proposition 5.1). We also provide a sharper $\tilde{O}(H^4 \sqrt{SAK})$ regret
 72 in the case when the missing rate is sufficiently small (Theorem 5.2).

73 To our best knowledge, these results present a first set of theories for RL with delayed and missing
 74 observations. Remarkably, our regret bounds nearly match the minimax-optimal regret of standard
 75 MDP in their dependence on S, A (noting that the target optimal policies are different in the two
 76 cases). It implies that RL with impaired observability are provably as efficient as RL with full
 77 observability (up to poly factors of H).

78 1.1 Related work

79 Efficient algorithms for learning in the standard setting of tabular MDPs without impaired observabil-
 80 ity has been extensively studied [Kearns and Singh, 2002, Brafman and Tennenholtz, 2002, Jaksch

81 et al., 2010, Dann and Brunskill, 2015, Azar et al., 2017, Agrawal and Jia, 2017, Jin et al., 2018,
 82 Dann et al., 2019, Zanette and Brunskill, 2019, Zhang et al., 2020, Domingues et al., 2021], where
 83 the minimax optimal regret is $\tilde{O}(\sqrt{H^3 SAK})$ [Azar et al., 2017, Domingues et al., 2021].

84 The delayed observation studied in this paper is related to delayed feedback in Howson et al. [2023],
 85 Yang et al. [2023], yet the setup is fundamentally different. In delayed feedback, an agent sends a
 86 policy to the environment for execution. The environment executes the policy on behalf of the agent
 87 for an episode, but the whole trajectory will be returned to the agent after some episodes. The policy
 88 executed by the environment is able to “see” instant state and reward. It is Markov and not played by
 89 the agent. Our setting concerns learning executable policies when delayed or missing states appear
 90 within an episode. The policy is no longer Markov and can only prescribe action based on history.
 91 Therefore, the algorithms and analyses for delayed feedback MDPs are not applicable to our settings.

92 Despite the distinct settings, there are existing fruitful results in efficiently learning MDPs or bandits
 93 with delayed feedback. Stochastic delayed feedback in bandits is studied in Agarwal and Duchi [2011],
 94 Dudik et al. [2011], Joulani et al. [2013], Vernade et al. [2017, 2020], Gael et al. [2020], Lancewicki
 95 et al. [2021]. In the more challenging setting of reinforcement learning, Howson et al. [2023]
 96 considers tabular MDPs and Yang et al. [2023] generalizes to MDPs with function approximation
 97 and multi-agent settings.

98 On the other hand, results analyzing MDPs with missing observations are limited in literature,
 99 although missing data is a commonly recognized issue in applications [García-Laencina et al., 2010,
 100 Jerez et al., 2010, Little et al., 2012, Emmanuel et al., 2021]. One notable result is Bouneffouf et al.
 101 [2020] for bandits with missing rewards.

102 **Notation:** For real numbers a, b , we denote $a \wedge b = \min\{a, b\}$. In episodic MDPs, we use the
 103 superscript k to denote the index of episodes, and the subscript h to denote the index of time.
 104 We denote $\mathbf{a}_{i:j} = \{a_i, \dots, a_j\}$ as the collection of actions from time i to j . For two probability
 105 distributions μ and ν , we denote their total variation distance as $\|\mu - \nu\|_{\text{TV}}$.

106 **MDP preliminary:** An episodic MDP is described by a tuple $(\mathcal{S}, \mathcal{A}, H, R, P)$, where \mathcal{S}, \mathcal{A} are
 107 state and action spaces, respectively, H is the horizon, $R = \{r_h\}_{h=1}^H$ is the reward function and
 108 $P = \{p_h\}_{h=1}^H$ is the transition probability. We primarily focus on tabular MDPs, where $S = |\mathcal{S}|$ and
 109 $A = |\mathcal{A}|$ are both finite. We also assume that the reward is uniformly bounded with $\|r_h\|_1 \leq 1$ for
 110 any h . An agent will interact with the environment for K episodes, hoping to find a good policy to
 111 maximize the cumulative reward. Within an episode, at the h -th step, the agent chooses an action
 112 based on the available information of the environment. After taking the action, the underlying
 113 environment produces a reward and transits to the next state. With full state observation, a policy π
 114 maps instant state s to an action a or an action distribution. Given such a policy π , the value function
 115 is $V_h^\pi(s_1) = \mathbb{E}^\pi \left[\sum_{h^o=h}^H r_h(s_{h^o}, a_{h^o}) \mid s_h \right]$, where \mathbb{E}^π is the policy induced expectation.

116 2 Problem formulation

117 In this work, we study MDPs with impaired observability. We focus on two practical settings: 1)
 118 delayed observations and 2) missing observations.

MDP with delayed observations In any episode, we denote $d_h \in \{0, 1, \dots\}$ as the observational
 delay of the state and reward at step h . That is, we receive s_h and r_h at time $h + d_h$. The delay time
 d_h can be dependent on the state s_h and action a_h at time h . To facilitate analysis, we denote the inter-
 arrival time between the arrival of observations for step h and $h + 1$ as $\Delta_h = d_{h+1} - d_h$. With delays,
 at time h , the nearest observable state is denoted as s_{t_h} , where $t_h = \arg\max \{I : \sum_{i=0}^I \Delta_i \leq h\}$.
 Then the executable policy class

$$\Pi_e = \{\pi_h(\cdot \mid s_{t_h}, \mathbf{a}_{t_h:h-1}) \text{ for } h = 1, \dots, H\}$$

119 chooses actions depending on the nearest visible state and history actions. We impose the following
 120 assumption on the interarrival time.

121 **Assumption 2.1 .** The interarrival time Δ_h takes value in $\{0, 1, \dots\}$. The distribution $\mathcal{D}_h(s_h, a_h)$ of
 122 Δ_h can depend on (s_h, a_h) , but is conditionally independent of the MDP transitions given (s_h, a_h) .

123 Assumption 2.1 does not impose any distributional assumption on Δ_h , but only requires that the
 124 delayed observations still arrive in order and at each time step, there is at most one new visible state

125 and reward pair $(\Delta_h \geq 0)$. A widely studied example of delays in literature is that the inter-arrival
 126 time is geometrically distributed [Winsten, 1959]. Then the observation sequence $\{h + d_h\}$ is known
 127 as a Bernoulli process, which is understood as the discretized version of a Poisson process.

128 Our delayed observation setting is newly proposed and substantially generalizes the Constant Delayed
 129 MDPs (CDMDPs) studied in Brooks and Leondes [1972], Bander and White III [1999], Katsikopoulos
 130 and Engelbrecht [2003], Walsh et al. [2007]. When $\Delta_h = 0$ being deterministic for all $h \geq 1$ and k ,
 131 our observation delay coincides with CDMDPs. In CDMDPs, a new past observation is guaranteed
 132 to arrive at each time step. However, our delayed model can result in no new observation at some
 133 time steps.

134 Observation delay leads to difficulty in planning, as the agent can only infer the current state and then
 135 choose an action. Therefore, the policy is naturally history dependent. We summarize the interaction
 protocol of the agent with the environment in Protocol 1. At the end of each episode, we can collect

Protocol 1 Interaction between the agent and the environment with delayed observations

- 1: **for** episode $k = 1, \dots, K$ **do**
 - 2: **for** time $h = 1, \dots, H$ **do**
 - 3: The agent observes a pair of new, if any, state and reward $(s_{t_h}^k, a_{t_h}^k)$. By memory, the agent
 also has access to past actions $\mathbf{a}_{t_h:h-1}^k$.
 - 4: The agent plays action a_h^k according to some executable policy $\pi_h^k \in \Pi_e$.
 - 5: The environment transits to next state $s_{h+1}^k \sim p_h(\cdot | s_h^k, a_h^k)$, which is unobservable to the
 agent. The environment also decides the delay at step $h + 1$ as $d_{h+1}^k = d_h^k + \Delta_h^k$ and t_{h+1}^k .
 - 6: **end for**
 - 7: The environment sends all unobserved pairs of state and reward as well as their corresponding
 delay time to the agent.
 - 8: **end for**
-

139 all the delayed observations, however, these observations are not used in planning. In reality, the
 138 agent can collect these observations by waiting after time H .

139 **MDP with missing observations** In addition to the stochastic delay in observations, we also
 140 consider randomly missing observations. In applications, an agent interacts with the environment
 141 through some communication channel. The communication channel is often imperfect and thus,
 142 observation can be lost during transmission. This type of missing is permanent and we describe in
 143 the following assumption.

144 **Assumption 2.2** . Any pair of observation (state and reward) is independently observable in the
 145 communication channel. The observation rate is λ_h depending on h , but independent of the MDP
 146 transitions. Moreover, there exists a constant λ_0 such that $\lambda_h \geq \lambda_0$ for any h . The agent will be
 147 informed when an observation is missing.

148 3 Construction of augmented MDPs

149 To tackle the limited observability, we expand the original state space and define an augmented
 150 MDP. It will serve as the basis for our subsequent theoretical analysis. For audience not interested in
 151 technical details, please feel free to skip this section.

152 3.1 Augmented MDP with expected reward

153 In the remainder of this section, we focus on the delayed observation case and defer the missing case
 154 to Section 5. Define $\tau_h = \{s_{t_h}, \mathbf{a}_{t_h:h-1}, \delta_{t_h}\}$ as the augmented state, where $\delta_{t_h} \in [0, \Delta_{t_h}]$ is the
 155 delayed steps after observing (s_{t_h}, r_{t_h}) . Let \mathcal{S}_{aug} denote the augmented state space of all possible
 156 τ 's. Then the original MDP with delayed observations can be reformulated into a state-augmented
 157 one $\text{MDP}_{\text{aug}} = (\mathcal{S}_{\text{aug}}, \mathcal{A}, H, R_{\text{aug}}, P_{\text{aug}})$. The reward is defined as

$$r_{h,\text{aug}}(\tau_h, a_h) = \mathbb{E}[r_h(s_h, a_h) | \tau_h, a_h],$$

158 which is the expected reward given the nearest past state s_{t_h} and history actions $\mathbf{a}_{t_h:h-1}$. We can define
 159 belief distribution $b_h(s | \tau_h) = P(s_h = s | \tau_h)$. Then $r_{h,\text{aug}}(\tau_h, a_h) = \mathbb{E}_{s \sim b_h(\cdot | \tau_h)}[r(s, a_h)]$. Belief
 160 distributions are widely adopted in partially observed MDPs [Ross et al., 2007, Poupart and Vlassis,
 161 2008]. We will frequently use the belief distribution to study the expressivity of Π_e in Section 4.2.

162 The transition probabilities P_{aug} are sparse. For any $\tau_h = \{s_{t_h}, \mathbf{a}_{t_h:h-1}, \delta_{t_h}\}$ and $\tau_{h+1} =$
 163 $\{s_{t_{h+1}}, \mathbf{a}_{t_{h+1}:h}, \delta_{t_{h+1}}\}$, we have

$p_{h,\text{aug}}(\tau_{h+1} \tau_h, a_h)$	Condition
$M_a(\tau_h, \tau_{h+1})\theta_{\text{delay}}(s_{t_h}, a_{t_h}, \delta_{t_h})p_h(s_{t_{h+1}} s_{t_h}, a_{t_h})$	if $\delta_{t_{h+1}} = 0$ and $t_{h+1} = t_h + 1$
$M_a(\tau_h, \tau_{h+1})(1 - \theta_{\text{delay}}(s_{t_h}, a_{t_h}, \delta_{t_h}))$	if $\delta_{t_{h+1}} = \delta_{t_h} + 1$ and $t_{h+1} = t_h$
0	otherwise

164 where $M_a(\tau_h, \tau_{h+1})$ indicates whether the rolling actions are matched, i.e.,

$$M_a(\tau_h, \tau_{h+1}) = \mathbb{1}\{\mathbf{a}_{t_h:h-1} = \mathbf{a}_{t_{h+1}:h-1}\},$$

165 and $\theta_{\text{delay}}(s_{t_h}, a_{t_h}, \delta_{t_h})$ is defined as

$$\theta_{\text{delay}}(s_{t_h}, a_{t_h}, \delta_{t_h}) = \mathbb{P}(\Delta_{t_h} = \delta_{t_h} | s_{t_h}, a_{t_h}, \delta_{t_h}) = \frac{\mathbb{P}(\Delta_{t_h} = \delta_{t_h} | s_{t_h}, a_{t_h})}{1 - \sum_{\delta < \delta_{t_h}} \mathbb{P}(\Delta_{t_h} = \delta | s_{t_h}, a_{t_h})}.$$

166 The factored form of $\theta_{\text{delay}}(s_{t_h}, a_{t_h}, \delta_{t_h})p_h(s_{t_{h+1}}|s_{t_h}, a_{t_h})$ follows from the conditional independence
 167 in Assumption 2.1. We define Q -functions and value functions as follows. For any τ_h, a_h and policy
 168 $\pi \in \Pi_e$, we have

$$Q_{h,\text{aug}}^\pi(\tau_h, a_h) = \mathbb{E}^\pi \left[\sum_{h^0=h}^H r_{h,\text{aug}}(\tau_{h^0}, a_{h^0}) \middle| \tau_h, a_h \right] \quad \text{and}$$

$$V_{h,\text{aug}}^\pi(\tau_h) = \langle Q_{h,\text{aug}}^\pi(\tau_h, \cdot), \pi_h(\cdot | \tau_h) \rangle.$$

169 We note that V_h^π is equivalent to $V_{h,\text{aug}}^\pi$ for the same executable policy $\pi \in \Pi_e$. We also denote $\mathcal{P}_{h,\text{aug}}$
 170 as the transition operator corresponding to P_{aug} . It can be checked that

$$Q_{h,\text{aug}}^\pi(\tau_h, a_h) = r_{h,\text{aug}}(\tau_h, a_h) + [\mathcal{P}_{h,\text{aug}} V_{h,\text{aug}}^\pi](\tau_h, a_h).$$

171 MDP_{aug} also appears in makes all the policies in Π_e executable and Markov. Meanwhile, the reward
 172 function keeps track of all the expected reward for H steps. Although the expanded state space
 173 \mathcal{S}_{aug} is much more complicated than the original state space \mathcal{S} , the sparse structures in the transition
 174 probabilities still allow an efficient exploration. We note that $p_{h,\text{aug}}$ only depends on the delay
 175 distribution and one-step Markov transitions. However, there is still one caveat for learning in MDP_{aug}
 176 – the reward function depends belief distributions, which involve multi-step transitions.

177 3.2 Augmented MDP with past reward

178 To tackle the aforementioned challenge, we further define $\widetilde{\text{MDP}}_{\text{aug}} = (\widetilde{\mathcal{S}}_{\text{aug}}, \mathcal{A}, \widetilde{H}, \widetilde{R}_{\text{aug}}, \widetilde{P}_{\text{aug}})$ that
 179 shares the optimal policy in MDP_{aug} with an elonged horizon $\widetilde{H} = 2H$. The state space \mathcal{S}_{aug} consists
 180 of any $\tau_h = \{s_{t_h}, \mathbf{a}_{t_h:h \wedge H}, \delta_{t_h}\}$. Comparing to \mathcal{S}_{aug} , we cut off the action at horizon H , since a_h
 181 for $h > H$ has no influence on the state and reward in time $[0, H]$. The reward function is defined as

$$\widetilde{r}_{h,\text{aug}}(\tau_h, a_h) = r_h(s_{t_h}, a_{t_h}) \mathbb{1}\{\delta_{t_h} = 0\} \mathbb{1}\{t_h \in \{1, \dots, H\}\}.$$

182 By definition, $\widetilde{r}_{h,\text{aug}}(\tau_h, a_h)$ is a past reward. More importantly, $\widetilde{r}_{h,\text{aug}}(\tau_h, a_h)$ zeros out rewards
 183 outside the original horizon H . Meanwhile, between the arrival of two consecutive state observations,
 184 the reward only counts once. Lastly, the transition probabilities are

$\widetilde{p}_{h,\text{aug}}(\tau_{h+1} \tau_h, a_h)$	Condition
$M_a(\tau_h, \tau_{h+1})\theta_{\text{delay}}(s_{t_h}, a_{t_h}, \delta_{t_h})p_h(s_{t_{h+1}} s_{t_h}, a_{t_h})$	if $\delta_{t_{h+1}} = 0, t_{h+1} = t_h + 1$ and $h < H$
$M_a(\tau_h, \tau_{h+1})(1 - \theta_{\text{delay}}(s_{t_h}, a_{t_h}, \delta_{t_h}))$	if $\delta_{t_{h+1}} = \delta_{t_h} + 1, t_{h+1} = t_h$ and $h < H$
$M_a(\tau_h, \tau_{h+1})p_h(s_{t_{h+1}} s_{t_h}, a_{t_h})$	if $\delta_{t_{h+1}} = 0, t_{h+1} = t_h + 1$ and $h > H$
0	otherwise

185 We interpret the transitions as follows. When $h \leq H$, the transition is the same as MDP_{aug} . When
 186 $h > H$, we simply wait for unobserved states and rewards to come. As mentioned, actions taken
 187 beyond time H are irrelevant. We build an equivalence in the expected values of MDP_{aug} and $\widetilde{\text{MDP}}_{\text{aug}}$.

188 **Proposition 3.1.** Let MDP_{aug} and $\widetilde{\text{MDP}}_{\text{aug}}$ be defined as in the previous paragraphs. Then for any
 189 initial state τ_1 and any policy $\pi = \{\pi_h\}_{h=1}^H \in \Pi_e$, it holds,

$$\mathbb{E}^\pi \left[\sum_{h=1}^H r_{h,\text{aug}}(\tau_h, a_h) \middle| \tau_1 \right] = \mathbb{E}^\pi \left[\sum_{h=1}^{\widetilde{H}} \widetilde{r}_{h,\text{aug}}(\tau_h, a_h) \middle| \tau_1 \right],$$

190 where in the right-hand side, the policy for steps $H + 1$ to \widetilde{H} is arbitrary.

191 The proof is provided in Appendix A.1. Proposition 3.1 implies that learning in MDP_{aug} until time H
 192 is equivalent to that in $\widetilde{\text{MDP}}_{\text{aug}}$ for \widetilde{H} steps.

193 4 RL with delayed observations and regret bound

194 In this section, we provide regret analysis of learning in MDPs with stochastic delays. For the sake of
 195 simplicity, we assume the reward is known, however, extension to unknown reward causes no real
 196 difficulty. Motivated by the augmented MDP reformulation, we introduce our learning algorithm
 197 in Algorithm 2. In Line 5, unobserved states and rewards are returned to the agent as described
 198 in Protocol 1. Using the data set, we construct bonus functions compensating the uncertainty in
 199 *one-step* transitions of the original MDP. This largely sharpens the confidence region, yet still ensures
 200 a valid optimism. We emphasize that in Line 9, we are planning on $\widetilde{\text{MDP}}_{\text{aug}}$ involving the augmented
 201 transitions and expanded states of $\tau \in \widetilde{\mathcal{S}}_{\text{aug}}$. Only in this way, we can obtain an executable policy in
 delayed MDPs. The planning complexity is SA^H though.

Algorithm 2 Policy learning for delayed MDPs using $\widetilde{\text{MDP}}_{\text{aug}}$

- 1: **Input:** Original horizon H , extended horizon \widetilde{H} , policy class Π_e , failure probability γ .
- 2: **Init:** $V_{\widetilde{H}+1}(\tau) = 0$ and $Q_{\widetilde{H}}(\tau, a) = H$ for any τ and a , data set $\mathcal{D}^0 = \emptyset$, initial policy π^0 .
- 3: **for** episode $k = 1, \dots, K$ **do**
- 4: Execute policy π^{k-1} for \widetilde{H} steps.
- 5: After the episode ends, collect data $\mathcal{D}^k = \mathcal{D}^{k-1} \cup \{(s_h^k, a_h^k, r_h^k, \Delta_h^k)\}_{h=1}^H$.
- 6: On data set \mathcal{D}^k , compute counting numbers $N_h^k(s_h, a_h)$, $N_h^k(s_h, a_h, s_{h+1})$ and
 $N_h^k(s_h, a_h, \delta_h)$.
- 7: Estimate transition probabilities and delay distributions via

$$\widehat{p}_h^k(s_{h+1}|s_h, a_h) = \frac{N_h^k(s_h, a_h, s_{h+1})}{N_h^k(s_h, a_h)} \quad \text{and} \quad \widehat{\theta}_{\text{delay}}^k(s_h, a_h, \delta_h) = \frac{N_h^k(s_h, a_h, \delta_h)}{\sum_{\delta} N_h^k(s_h, a_h, \delta)}.$$

Then estimators of $\widetilde{p}_{h,\text{aug}}$ in $\widetilde{\text{MDP}}_{\text{aug}}$ is computed using \widehat{p}_h^k and $\widehat{\theta}_{\text{delay}}^k$.

- 8: Set bonus function as

$$b_h^k(\tau_h, a_h) = cH \left(\sqrt{\frac{H\iota}{N_{t_h}^k(s_{t_h}, a_h, \delta_{t_h})}} + \sqrt{\frac{H\iota}{N_{t_h}^k(s_{t_h}, a_{t_h})}} \right)$$

for $\iota = \log \frac{SAKH}{\gamma}$ and c sufficiently large.

- 9: Run optimistic value iteration in $\widetilde{\text{MDP}}_{\text{aug}}$ for \widetilde{H} steps and obtain $\pi^k \in \Pi_e$.
 - 10: **end for**
 - 11: **Return:** Learned policy $\pi_{1:H}^k$ for $k = 1, \dots, K$.
-

202

203 4.1 Regret bound

204 We define regret in delayed MDP as

$$\text{Regret}(K) = \sum_{k=1}^K \max_{\pi \in \Pi_e} V_1^{\pi}(s_1^k) - \sum_{k=1}^K V_1^{\pi^k}(s_1^k),$$

205 where V_1^{π} is the value function of the original MDP. Although the regret here is defined on the original
 206 MDP, it is equivalent to the regret of the same policy on MDP_{aug} and further $\widetilde{\text{MDP}}_{\text{aug}}$ by Proposition 3.1.
 207 Note that we are comparing with the best executable policy. The performance degradation caused by
 208 observation delay is discussed in Section 4.2. The following theorem bounds the regret.

209 **Theorem 4.1** (Regret bound for Delayed MDP). Suppose Assumption 2.1 holds. Let $\gamma \in (0, 1)$ be
 210 any failure probability. With probability $1 - \gamma$, the regret of Algorithm 2 satisfies

$$\text{Regret}(K) \leq c \left(H^4 \sqrt{SAK\iota} + H^4 S^2 A \iota^2 \right),$$

211 where $\iota = \log \frac{SAHK}{\gamma}$ and c is some constant.

212 The proof is provided in Appendix B.1. We discuss several implications.

213 **Sharp dependence on S and A** Theorem 4.1 has a sharp dependence on S and A , although the
 214 expanded state space \tilde{S}_{aug} has a cardinality bounded by SA^H . Naively learning and planning in
 215 $\widehat{\text{MDP}}_{\text{aug}}$ would suffer from the exponential enlargement of A^H . However, we identify the sparse
 216 structures in the transition probabilities. As can be seen, $\tilde{p}_{h,\text{aug}}$ only involves one-step transitions in
 217 the original MDP and some conditionally independent delay distributions. Such structures lead to a
 218 rather easy estimation of $\tilde{p}_{h,\text{aug}}$, which can be constructed from the estimators of one-step transitions
 219 in the original MDP. Meanwhile, the sparse structures make exploration in $\widehat{\text{MDP}}_{\text{aug}}$ efficient, due to
 220 many unreachable states.

221 **Effect of the delay distribution and delay length** Theorem 4.1 holds for arbitrary conditionally
 222 independent delay distributions, even include heavy-tailed distributions. Our regret bound encodes
 223 the influence of delay by paying extra H factors. The reason to this is that if the delay is larger than
 224 H , then the corresponding state will only be observed after an episode ends and won't be used in
 225 planning. Therefore, we can truncate the delay at H , regardless of its tail distributions.

226 4.2 Performance degradation of policy class Π_e

227 This section devotes to quantify the performance degradation caused by delayed observations. In
 228 particular, we bound the value difference between the best executable policy and the best Markov
 229 policy in a no delay environment. Recall that V_1 is the value function of the original MDP. We denote

$$\pi_{\text{nodelay}} = \operatorname{argmax}_{\pi} V_1^{\pi}(s_1) \quad \text{and} \quad \pi_{\text{delay}} = \operatorname{argmax}_{\pi \in \Pi_e} V_1^{\pi}(s_1)$$

230 as the best vanilla optimal policy and executable policy, respectively. The values achieved by π_{nodelay}
 231 and π_{delay} are denoted as $V_{1,\text{nodelay}}(s_1)$ and $V_{1,\text{delay}}(s_1)$, respectively. The gap between $V_{1,\text{nodelay}}$
 232 and $V_{1,\text{delay}}$ quantifies the performance degradation, which is denoted as $\text{gap}(s_1) = V_{1,\text{nodelay}}(s_1) -$
 233 $V_{1,\text{delay}}(s_1)$. We bound gap in Proposition B.2 in Appendix due to space limit.

234 In a nutshell, we show that the performance degradation gap is highly relevant to the belief distribution
 235 $b_h(\cdot|\tau)$. When $b_h(\cdot|\tau)$ is evenly spread, meaning that the entropy of b_h is high and inferring the
 236 current unseen state is difficult, we potentially suffer from a large gap. On the contrary, when $b_h(\cdot|\tau)$
 237 is nearly deterministic, the performance degradation is small. In the special case of deterministic
 238 transitions, we have $\text{gap} = 0$.

239 4.3 The (mysterious) effect of delay on the optimal value

240 To further understand the effect of the delay on the optimal value, we provide the following dichotomy.
 241 On the one hand, we show that there exists an MDP instance, such that a constant delay of d steps
 242 does not hurt the performance. On the other hand, in the same MDP instance, a constant delay of
 243 $d + 1$ steps suffers from a constant performance drop.

244 **Proposition 4.2.** Consider constant delayed MDPs. Fix a positive integer $d < H$. Then there exists
 245 an MDP instance such that the following two items hold simultaneously.

- 246 • When delay is d , it holds $\frac{1}{K} \sum_{k=1}^K \text{gap}(s_1^k) = 0$.
- 247 • When delay is $d + 1$, it holds $\frac{1}{K} \sum_{k=1}^K \text{gap}(s_1^k) \geq \frac{1}{2} - \sqrt{\frac{1}{2K} \log \frac{1}{\gamma}}$, with probability $1 - \gamma$.

248 The proof is provided in Appendix B.3. We remark that Proposition 4.2 says that observation delay
 249 can be dangerous, even with the slightest possible number of steps. The idea behind Proposition 4.2
 250 is consistent with the analysis on gap. In particular, we construct an MDP instance demonstrated
 251 in Figure 2. The reward vanishes at all times but $d + 1$. When delay is d , the initial state s_1 is
 252 revealed and the policy can choose the best action to receive a reward. When delay is $d + 1$, however,
 253 there is always a $1/2$ probability of missing the best action for any policy, which leads to a constant
 254 performance degradation.

255 5 RL with missing observations and regret analysis

256 We now switch our study to MDPs with missing observations. In such an environment, executable
 257 policies share the same structures as delayed MDPs, where an action is taken based on available
 258 history information. Compared to delayed observations, learning with missing observations is
 259 more challenging. Since unobserved states and rewards are never revealed, we are suffering from
 260 information loss. Besides, we will frequently deal with multi-step transitions, due to missing
 261 observations between two consecutive visible states.

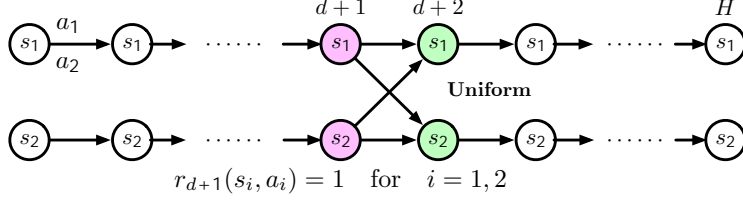


Figure 2: MDP instance on two states with two actions. The transition is lazy until time d . Then the transition is uniform regardless of actions for time $d + 1$. Reward is nonzero only at time $d + 1$. This is an example with a delay of length d causes no degradation and a delay of $d + 1$ causes a constant performance degradation.

262 5.1 Optimistic planning with missing observations

263 Despite the difficulty, we present here algorithms that are efficient in learning and planning for MDPs
 264 with missing observations. We begin with an optimistic planning algorithm in Algorithm 3. To unify
 265 the notation, we denote $s_h^k = \emptyset$ and $r_h^k = \emptyset$ as missing the observation.

Algorithm 3 Optimistic planning for MDPs with missing observations

- 1: **Input:** Horizon H , observable rate λ_h .
- 2: **Init:** $\mathcal{B}^0 = \Theta$ to be all possible tabular MDPs, data set $\mathcal{D}^0 = \emptyset$.
- 3: **for** episode $k = 1, \dots, K$ **do**
- 4: Set policy $\pi^k = \operatorname{argmax}_{\pi} \max_{\theta \in \Theta} V_{1,\theta}^{\pi}(s_1^k)$.
- 5: Play policy π^k and collect data $\mathcal{D}^{k-1} \cup \{(s_h^k, a_h^k, r_h^k)\}_{h=1}^H$.
- 6: Compute counting number $N_h^k(s, a) = \sum_{j=1}^k \mathbf{1}\{s_j^k = s, a_j^k = a, s_{j+1}^k \neq \emptyset\}$.
- 7: Update confidence set

$$\mathcal{B}^k = \left\{ \theta : \|\hat{p}_h^k(\cdot|s, a) - p_h^\theta(\cdot|s, a)\|_{\text{TV}} \leq c \sqrt{\frac{S\iota}{N_h^k(s, a)}} \text{ for all } (h, s, a) \right\} \cap \mathcal{B}^{k-1},$$

where $\hat{p}_h^k(s^\emptyset|s, a) = \frac{N_h^k(s, a, s^\emptyset)}{N_h^k(s, a)}$ and c is some constant.

8: **end for**

266 The majority of the algorithm resembles the typical optimistic planning [Jaksch et al., 2010] but with
 267 some notable differences. In Line 4, the value function $V_{1,\theta}$ is for the original MDP with transition
 268 probabilities parameterized by θ . Different from the typical optimistic planning, the underlying MDP
 269 here obeys the stochastic observable model in Assumption 2.2. Therefore, the value $V_{1,\theta}$ is the sum
 270 of all possible values under missing observations. When counting $N_h^k(s, a)$ in Line 6, we exclude
 271 data tuples missing the next state, which inevitably slows down the learning curve. Nonetheless, the
 272 effect of missing only contributes as a scaling factor in the regret.

273 **Proposition 5.1.** Suppose Assumption 2.2 holds with λ_h known. Given a failure probability γ , with
 274 probability $1 - \gamma$, the regret of Algorithm 4 satisfies

$$\text{Regret}(K) \leq c \left(\frac{1}{-\log(1 - \lambda_0^2)} \sqrt{H^3 S^2 A K \iota^3} + \sqrt{H^4 K \iota} \right),$$

275 where $\iota = \log \frac{SAHK}{\gamma}$ and c is some constant.

276 The proof is provided in Appendix C.1. Proposition 5.1 is optimal in the K dependence and
 277 achieves an $S^2 A$ dependence on the complexity of the underlying MDP. In the extreme case of
 278 $\lambda_0 \approx 0$, which implies that every state and reward are hardly observable, we have $\text{Regret}(K) =$
 279 $\tilde{O}\left(\frac{1}{\lambda_0^2} \sqrt{H^3 S^2 A K}\right)$. Here λ_0^2 is the probability of observing two consecutive states for estimating
 280 the transition probabilities. Proposition 5.1 requires knowledge of observable rate λ_h . This is not a
 281 restrictive condition, as estimating λ_h from Bernoulli random variables is much easier than estimating
 282 transition probabilities.

283 **5.2 Model-based planning using augmented MDPs**

284 Proposition 5.1 has a lenient dependence on the missing rate $1 - \lambda_0^2$, nonetheless, is not sharp on the
 285 dependence of S . We next show that the augmented MDP approach is effective to tackle missing
 286 observations, when the observable rate satisfies additional conditions. Specifically, we assume that
 287 the observable rate λ_h is independent of (s, a) . We utilize the MDP_{aug} reformulation, except that we
 288 redefine the transition probabilities as

$$p_{h,\text{aug}}(\tau_{h+1} | \tau_h, a_h) = \begin{cases} \lambda_h p_h(s_{h+1} | s_{t_h}, \mathbf{a}_{t_h:h}) & \text{if } t_{h+1} = h + 1 \\ M_a(\tau_{h+1}, \tau_h)(1 - \lambda_h) & \text{if } t_{h+1} = t_h \\ 0 & \text{otherwise} \end{cases}.$$

289 The first case in $p_{h,\text{aug}}$ corresponds to receiving the state observation at time $h + 1$. In contrast to the
 290 delayed MDPs, the transition probabilities here potentially rely on multi-step transitions in the original
 291 MDP. The second case of the transition corresponds to missing the observation. We summarize the
 292 policy learning procedure in Algorithm 4 in Appendix C.2, which is similar to Algorithm 2, but with
 293 a new bonus function. The following theorem shows that Algorithm 4 is asymptotically efficient
 294 when the observable rate is relatively high.

295 **Theorem 5.2.** Suppose Assumption 2.2 holds with $\lambda_0 \geq 1 - A^{-(1+v)}$ for some positive constant v .
 296 Given a failure probability γ , with probability $1 - \gamma$, the regret of Algorithm 4 satisfies

$$\text{Regret}(K) \leq c \left(H^4 \sqrt{SAK} \iota^3 + S^2 \sqrt{H^9 K^{\frac{1}{1+v}} \iota^6} \right),$$

297 where $\iota = \log \frac{SAHK}{\gamma}$ and c is some constant.

298 The proof is provided in Appendix C.2. Some remarks are in order.

299 **SA rate when K is large** When the number of episodes $K \geq S^{3(1+v)/v}$, the first term
 300 $H^4 \sqrt{SAK} \iota^3$ in the regret bound dominates and attains a sharp dependence on S and A . How-
 301 ever, when the number of episodes are limited, the regret bound has a worse dependence on the state
 302 space size S . We also observe that as the missing rate λ becomes small (equivalently, v becomes
 303 large), the regret is close to $\tilde{O}(H^4 \sqrt{SAK} \iota^3)$.

304 **Observable rate smaller than $1 - 1/A$** Theorem 5.2 holds for an observable rate $\lambda_0 > 1 - 1/A$.
 305 The intuition behind is that to fully explore all the actions when a state observation is missing takes
 306 A trials. Therefore, in expectation, we will encounter a missing observation at least every A episodes
 307 as long as $\lambda_0 > 1 - 1/A$. Nonetheless, when $\lambda_0 \leq 1 - 1/A$, the regret bound remains curiously
 308 underexplored. We conjecture that $\lambda_0 = 1 - 1/A$ is a critical point distinguishes unique strategies
 309 for learning and planning in MDPs with missing observations. A detailed analysis goes beyond the
 310 scope of the current paper.

311 **Proof sketch** The proof of Theorem 5.2 adapts the analysis of model-based UCBVI algorithms
 312 [Azar et al., 2017]. Let m denote the maximal length of consecutive missing observations. We denote
 313 \mathcal{E}_m as the event when the maximal length of consecutive missing is less than m . On event \mathcal{E}_m , a naïve
 314 analysis leads to a $\tilde{O}(\sqrt{\text{poly}(H)SA^{m+1}K})$ regret, in observation to the size of the expanded state
 315 space \mathcal{S}_{aug} . However, our analysis circumvents the A^m dependence by exploiting the occurrence of
 316 consecutive missing observations is rare (Lemma C.3). On the complement of event, the regret is
 317 bounded by $KH(1 - \mathbb{P}(\mathcal{E}_m))$. Summing up the two parts and choosing a proper m yield our result.

318 **6 Conclusion and limitation**

319 In this paper, we have studied learning and planning in impaired observability MDPs. We focus
 320 on MDPs with delayed and missing observations. Specifically, for delayed observations, we have
 321 shown an efficient $\tilde{O}(H^4 \sqrt{SAK})$ regret. For missing observations, we have provided an optimistic
 322 planning algorithm achieving an $\tilde{O}(\sqrt{H^3 S^2 AK})$ regret. If the missing rate is relatively small, we
 323 have shown an efficient $\tilde{O}(H^4 \sqrt{SAK})$ regret bound. Further, we have characterized the performance
 324 drop caused by impaired observability compared to full observability. A limitation of the current
 325 study is that the planning complexity in augmented MDPs is high with an exponential dependence on
 326 the size of the action space. Sharpening such a dependence is left as a future direction.

327 **References**

- 328 Alekh Agarwal and John C Duchi. Distributed delayed stochastic optimization. *Advances in Neural*
329 *Information Processing Systems*, 24, 2011.
- 330 Mridul Agarwal and Vaneet Aggarwal. Blind decision making: Reinforcement learning with delayed
331 observations. *Pattern Recognition Letters*, 150:176–182, 2021.
- 332 Shipra Agrawal and Randy Jia. Optimistic posterior sampling for reinforcement learning: worst-case
333 regret bounds. *Advances in Neural Information Processing Systems*, 30, 2017.
- 334 Mohammad Gheshlaghi Azar, Ian Osband, and Rémi Munos. Minimax regret bounds for rein-
335 forcement learning. In *Proceedings of the International Conference on Machine Learning*, pages
336 263–272. PMLR, 2017.
- 337 James L Bander and Chelsea C White III. Markov decision processes with noise-corrupted and
338 delayed state observations. *Journal of the Operational Research Society*, 50(6):660–668, 1999.
- 339 Dimitri Bertsekas. *Dynamic Programming and Optimal Control: Volume I*, volume 1. Athena
340 Scientific, 2012.
- 341 Djallel Bouneffouf, Sohini Upadhyay, and Yasaman Khazaeni. Contextual bandit with missing
342 rewards. arXiv preprint arXiv:2007.06368, 2020.
- 343 Ronen I Brafman and Moshe Tennenholtz. R-max-a general polynomial time algorithm for near-
344 optimal reinforcement learning. *Journal of Machine Learning Research*, 3(Oct):213–231, 2002.
- 345 DM Brooks and Cornelius T Leondes. Markov decision processes with state-information lag.
346 *Operations Research*, 20(4):904–907, 1972.
- 347 Fan Chen, Yu Bai, and Song Mei. Partially observable RL with B-stability: Unified structural
348 condition and sharp sample-efficient algorithms. arXiv preprint arXiv:2209.14990, 2022.
- 349 Fan Chen, Huan Wang, Caiming Xiong, Song Mei, and Yu Bai. Lower bounds for learning in
350 revealing POMDPs. arXiv preprint arXiv:2302.01333, 2023.
- 351 Christoph Dann and Emma Brunskill. Sample complexity of episodic fixed-horizon reinforcement
352 learning. *Advances in Neural Information Processing Systems*, 28, 2015.
- 353 Christoph Dann, Lihong Li, Wei Wei, and Emma Brunskill. Policy certificates: Towards accountable
354 reinforcement learning. In *Proceedings of the International Conference on Machine Learning*,
355 pages 1507–1516. PMLR, 2019.
- 356 Omar Darwiche Domingues, Pierre Ménard, Emilie Kaufmann, and Michal Valko. Episodic rein-
357 forcement learning in finite mdps: Minimax lower bounds revisited. In *Algorithmic Learning*
358 *Theory*, pages 578–598. PMLR, 2021.
- 359 Miroslav Dudik, Daniel Hsu, Satyen Kale, Nikos Karampatziakis, John Langford, Lev Reyzin, and
360 Tong Zhang. Efficient optimal learning for contextual bandits. arXiv preprint arXiv:1106.2369,
361 2011.
- 362 Tlamele Emmanuel, Thabiso Maupong, Dimane Mpoeleng, Thabo Semong, Banyatsang Mphago,
363 and Oteng Tabona. A survey on missing data in machine learning. *Journal of Big Data*, 8(1):1–37,
364 2021.
- 365 Manegueu Anne Gael, Claire Vernade, Alexandra Carpentier, and Michal Valko. Stochastic bandits
366 with arm-dependent delays. In *Proceedings of the International Conference on Machine Learning*,
367 pages 3348–3356. PMLR, 2020.
- 368 Pedro J García-Laencina, José-Luis Sancho-Gómez, and Aníbal R Figueiras-Vidal. Pattern classifica-
369 tion with missing data: a review. *Neural Computing and Applications*, 19:263–282, 2010.
- 370 Benjamin Howson, Ciara Pike-Burke, and Sarah Filippi. Delayed feedback in generalised linear
371 bandits revisited. In *International Conference on Artificial Intelligence and Statistics*, pages
372 6095–6119. PMLR, 2023.

- 373 Thomas Jaksch, Ronald Ortner, and Peter Auer. Near-optimal regret bounds for reinforcement
374 learning. *Journal of Machine Learning Research*, 11(51):1563–1600, 2010. URL [http://jmlr.](http://jmlr.org/papers/v11/jaksch10a.html)
375 [org/papers/v11/jaksch10a.html](http://jmlr.org/papers/v11/jaksch10a.html).
- 376 José M Jerez, Ignacio Molina, Pedro J García-Laencina, Emilio Alba, Nuria Ribelles, Miguel Martín,
377 and Leonardo Franco. Missing data imputation using statistical and machine learning methods in a
378 real breast cancer problem. *Artificial Intelligence in Medicine*, 50(2):105–115, 2010.
- 379 Chi Jin, Zeyuan Allen-Zhu, Sebastien Bubeck, and Michael I Jordan. Is Q-learning provably efficient?
380 *Advances in Neural Information Processing Systems*, 31, 2018.
- 381 Chi Jin, Sham Kakade, Akshay Krishnamurthy, and Qinghua Liu. Sample-efficient reinforcement
382 learning of undercomplete POMDPs. *Advances in Neural Information Processing Systems*, 33:
383 18530–18539, 2020.
- 384 Pooria Joulani, Andras Gyorgy, and Csaba Szepesvári. Online learning under delayed feedback. In
385 *Proceedings of the International Conference on Machine Learning*, pages 1453–1461. PMLR,
386 2013.
- 387 Konstantinos V Katsikopoulos and Sascha E Engelbrecht. Markov decision processes with delays
388 and asynchronous cost collection. *IEEE Transactions on Automatic Control*, 48(4):568–574, 2003.
- 389 Michael Kearns and Satinder Singh. Near-optimal reinforcement learning in polynomial time.
390 *Machine Learning*, 49(2):209–232, 2002.
- 391 Vikram Krishnamurthy. *Partially Observed Markov Decision Processes*. Cambridge University Press,
392 2016.
- 393 Tal Lancelwicki, Shahar Segal, Tomer Koren, and Yishay Mansour. Stochastic multi-armed bandits
394 with unrestricted delay distributions. In *Proceedings of the International Conference on Machine*
395 *Learning*, pages 5969–5978. PMLR, 2021.
- 396 Tor Lattimore and Csaba Szepesvári. *Bandit Algorithms*. Cambridge University Press, 2020.
- 397 Roderick J Little, Ralph D’Agostino, Michael L Cohen, Kay Dickersin, Scott S Emerson, John T
398 Farrar, Constantine Frangakis, Joseph W Hogan, Geert Molenberghs, Susan A Murphy, et al. The
399 prevention and treatment of missing data in clinical trials. *New England Journal of Medicine*, 367
400 (14):1355–1360, 2012.
- 401 Qinghua Liu, Praneeth Netrapalli, Csaba Szepesvari, and Chi Jin. Optimistic MLE—A generic
402 model-based algorithm for partially observable sequential decision making. arXiv preprint
403 arXiv:2209.14997, 2022.
- 404 Shichao Liu, Xiaoyu Wang, and Peter Xiaoping Liu. Impact of communication delays on secondary
405 frequency control in an islanded microgrid. *IEEE Transactions on Industrial Electronics*, 62(4):
406 2021–2031, 2014.
- 407 Daniel J Lizotte, Lacey Gunter, Eric Laber, and Susan A Murphy. Missing data and uncertainty in
408 batch reinforcement learning. In *Advances in Neural Information Processing Systems*, 2008.
- 409 Christos H Papadimitriou and John N Tsitsiklis. The complexity of markov decision processes.
410 *Mathematics of Operations Research*, 12(3):441–450, 1987.
- 411 Pascal Poupart and Nikos Vlassis. Model-based bayesian reinforcement learning in partially ob-
412 servable domains. In *Proceedings of the International Symposium on Artificial Intelligence and*
413 *Mathematics*, pages 1–2, 2008.
- 414 Stephane Ross, Brahim Chaib-draa, and Joelle Pineau. Bayes-adaptive pomdps. *Advances in Neural*
415 *Information Processing Systems*, 20, 2007.
- 416 Richard D Smallwood and Edward J Sondik. The optimal control of partially observable markov
417 processes over a finite horizon. *Operations Research*, 21(5):1071–1088, 1973.
- 418 Richard S Sutton and Andrew G Barto. *Reinforcement Learning: An Introduction*. MIT Press, 2018.

- 419 Masatoshi Uehara, Ayush Sekhari, Jason D Lee, Nathan Kallus, and Wen Sun. Provably efficient
420 reinforcement learning in partially observable dynamical systems. arXiv preprint arXiv:2206.12020,
421 2022.
- 422 Claire Vernade, Olivier Cappé, and Vianney Perchet. Stochastic bandit models for delayed conversions.
423 arXiv preprint arXiv:1706.09186, 2017.
- 424 Claire Vernade, Alexandra Carpentier, Tor Lattimore, Giovanni Zappella, Beyza Ermis, and Michael
425 Brueckner. Linear bandits with stochastic delayed feedback. In *Proceedings of the International
426 Conference on Machine Learning*, pages 9712–9721. PMLR, 2020.
- 427 Martin J Wainwright. *High-dimensional Statistics: A Non-asymptotic Viewpoint*, volume 48. Cam-
428 bridge University Press, 2019.
- 429 Thomas J Walsh, Ali Nouri, Lihong Li, and Michael L Littman. Planning and learning in environments
430 with delayed feedback. In *Proceedings of Machine Learning: ECML 2007: 18th European
431 Conference on Machine Learning*, Warsaw, Poland, September 17-21, 2007, pages 442–453.
432 Springer, 2007.
- 433 CB Winsten. Geometric distributions in the theory of queues. *Journal of the Royal Statistical Society:
434 Series B (Methodological)*, 21(1):1–22, 1959.
- 435 Yunchang Yang, Han Zhong, Tianhao Wu, Bin Liu, Liwei Wang, and Simon S Du. A reduction-based
436 framework for sequential decision making with delayed feedback. arXiv preprint arXiv:2302.01477,
437 2023.
- 438 Andrea Zanette and Emma Brunskill. Tighter problem-dependent regret bounds in reinforcement
439 learning without domain knowledge using value function bounds. In *Proceedings of the Interna-
440 tional Conference on Machine Learning*, pages 7304–7312. PMLR, 2019.
- 441 Wenhao Zhan, Masatoshi Uehara, Wen Sun, and Jason D Lee. PAC reinforcement learning for
442 predictive state representations. arXiv preprint arXiv:2207.05738, 2022.
- 443 Zihan Zhang, Yuan Zhou, and Xiangyang Ji. Almost optimal model-free reinforcement learning
444 via reference-advantage decomposition. *Advances in Neural Information Processing Systems*, 33:
445 15198–15207, 2020.
- 446 Han Zhong, Wei Xiong, Sirui Zheng, Liwei Wang, Zhaoran Wang, Zhuoran Yang, and Tong Zhang.
447 A posterior sampling framework for interactive decision making. arXiv preprint arXiv:2211.01962,
448 2022.

449 **A Omitted proof in Section 3**

450 **A.1 Proof of Proposition 3.1**

451 *Proof of Proposition 3.1.* Consider an arbitrary fixed inter-arrival pattern $\Delta_0, \Delta_1, \dots, \Delta_{H-1}$. We
 452 show that the expected accumulated rewards under this inter-arrival pattern are identical for MDP_{aug}
 453 and $\widetilde{\text{MDP}}_{\text{aug}}$. In $\widetilde{\text{MDP}}_{\text{aug}}$, we have

$$\begin{aligned} & \mathbb{E}^\pi \left[\sum_{h=1}^{\tilde{H}} \tilde{r}_{h,\text{aug}}(\tau_h, a_h) \mid \tau_1, \Delta_0, \dots, \Delta_{H-1} \right] \\ \stackrel{(i)}{=} & \mathbb{E}^\pi \left[\sum_{h=1}^{\tilde{H}} \tilde{r}_{t_h,\text{aug}}(s_{t_h}, a_{t_h}) \mathbf{1}\{\delta_{t_h} = 0\} \mathbf{1}\{t_h \in \{1, \dots, H\}\} \mid \tau_1, \Delta_0, \dots, \Delta_{H-1} \right] \\ \stackrel{(ii)}{=} & \mathbb{E}^\pi \left[\sum_{h=1}^H r(s_h, a_h) \mid \tau_1, \Delta_0, \dots, \Delta_{H-1} \right] \\ = & \mathbb{E}^\pi \left[\sum_{h=1}^H r_{h,\text{aug}}(\tau_h, a_h) \mid \tau_1, \Delta_0, \dots, \Delta_{H-1} \right], \end{aligned}$$

454 where equality (i) invokes the definition of $\tilde{r}_{h,\text{aug}}$ and equality (ii) eliminates zero reward terms.
 455 Now taking expectation over all possible inter-arrival patterns, we deduce

$$\mathbb{E}^\pi \left[\sum_{h=1}^{\tilde{H}} \tilde{r}_{\text{aug}}(\tau_h, a_h) \mid \tau_1 \right] = \mathbb{E}^\pi \left[\sum_{h=1}^H r_{h,\text{aug}}(s_h, a_h) \mid \tau_1 \right].$$

456 The proof is complete. □

457 **B Omitted proofs in Section 4**

458 **B.1 Proof of Theorem 4.1**

459 *Proof of Theorem 4.1.* We adapt the main steps from Azar et al. [2017] for proving the theorem. The
 460 proof consists of verifying a valid optimism and developing a regret analysis. We denote $\tilde{Q}_{h,\text{aug}}$ as
 461 the optimal Q -function for $\widetilde{\text{MDP}}_{\text{aug}}$. When analyzing the regret, we also denote $\tilde{Q}_{h,\text{aug}}^k$ as the optimal
 462 Q -function in the k -th episode.

463 **Valid optimism** To begin with, we verify that the choice of the bonus functions leads to a valid
 464 optimism in the following lemma.

465 **Lemma B.1.** Given any failure probability $\gamma < 1$, we set a bonus as

$$b_h^k(\tau_h, a_h) = c_A H \left(\sqrt{\frac{H\iota}{N_{t_h}(s_{t_h}, a_{t_h}, \delta_{t_h})}} + \sqrt{\frac{H\iota}{N_{t_h}(s_{t_h}, a_{t_h})}} \right),$$

466 where $\iota = \log\left(\frac{SAHK}{\gamma}\right)$ and c_A is a constant. Then with probability $1 - \gamma$, it holds

$$\tilde{Q}_{h,\text{aug}}^k(\tau_h, a_h) \geq \tilde{Q}_{h,\text{aug}}(\tau_h, a_h), \quad \tilde{V}_{h,\text{aug}}^k(\tau_h) \geq \tilde{V}_{h,\text{aug}}(\tau_h) \quad \text{for any } (k, h, \tau_h, a_h).$$

467 *Proof of Lemma B.1.* We compute the cardinality of the expanded state space $\tilde{\mathcal{S}}_{\text{aug}}$ as

$$|\tilde{\mathcal{S}}_{\text{aug}}| \stackrel{(i)}{=} \sum_{i=0}^H HSA^i = HS \frac{A^{H+1} - 1}{A - 1} \leq 2HSA^H.$$

468 For a fixed episode k , we show by backward induction that the assertion in Lemma B.1 holds. To ease
 469 the presentation, we omit all superscripts of k , all subscripts of “aug”, as well as the tilde $\tilde{\cdot}$ notation.

470 When $h = \tilde{H} + 1$, the base assertion holds immediately. Suppose the assertion is true for time $h + 1$.
 471 At time h , for any fixed (τ_h, a_h) , if $Q_h(\tau_h, a_h) = H$, the assertion holds true. Otherwise, we have

$$\begin{aligned} Q_h(\tau_h, a_h) - Q_h(\tau_h, a_h) &= [\hat{\mathcal{P}}_h V_{h+1}](\tau_h, a_h) - [\mathcal{P}_h V_{h+1}](\tau_h, a_h) + b_h^k(\tau_h, a_h) \\ &\geq \underbrace{\left([\hat{\mathcal{P}}_h - \mathcal{P}_h]V_{h+1}\right)}_{(A)}(\tau_h, a_h) + b_h^k(\tau_h, a_h). \end{aligned}$$

472 We show a lower bound on (A). If $h \geq H$, expanding the transition kernel \mathcal{P}_h leads to

$$\begin{aligned} (A) &= \sum_{\tau_{h+1}} V_{h+1}(\tau_{h+1})(\hat{p}_h(\tau_{h+1}|\tau_h, a_h) - p_h(\tau_{h+1}|\tau_h, a_h)) \\ &\stackrel{(i)}{=} \sum_{s_{t_h+1}} V_{h+1}(\tau_{h+1})(\hat{p}_h(s_{t_h+1}|s_{t_h}, a_{t_h}) - p_h(s_{t_h+1}|s_{t_h}, a_{t_h})) \\ &\stackrel{(ii)}{\geq} -c_{A,1}H \sqrt{\frac{H\iota}{N_{t_h}(s_{t_h}, a_{t_h})}}, \end{aligned}$$

473 where equality (i) requires τ_{h+1} to take s_{t_h+1} as the new state observation, and inequality (ii) follows
 474 from the Hoeffding's inequality (Lemma D.2) with a constant $c_{A,1}$. Note that the $H\iota$ term in the
 475 numerator comes from a union bound over $\tilde{\mathcal{S}}_{\text{aug}} \times \mathcal{A}$.

476 On the other hand, if $h < H$, expanding the transition kernel \mathcal{P}_h yields

$$\begin{aligned} (A) &= \sum_{\tau_{h+1}} V_{h+1}(\tau_{h+1})(\hat{p}_h(\tau_{h+1}|\tau_h, a_h) - p_h(\tau_{h+1}|\tau_h, a_h)) \\ &= \underbrace{\sum_{\tau_{h+1}} V_{h+1}(\tau_{h+1})(\hat{p}_h(\tau_{h+1}|\tau_h, a_h) - p_h(\tau_{h+1}|\tau_h, a_h)) \mathbf{1}\{\delta_{t_{h+1}} = 0\} \mathbf{1}\{t_{h+1} = t_h + 1\}}_{(A_1)} \\ &\quad + \underbrace{\sum_{\tau_{h+1}} V_{h+1}(\tau_{h+1})(\hat{p}_h(\tau_{h+1}|\tau_h, a_h) - p_h(\tau_{h+1}|\tau_h, a_h)) \mathbf{1}\{\delta_{t_{h+1}} = \delta_{t_h} + 1\} \mathbf{1}\{t_{h+1} = t_h\}}_{(A_2)}. \end{aligned}$$

477 Note that (A₁) accounts for receiving a new state observation in τ_{h+1} , and (A₂) accounts for no new
 478 state observation. We tackle these two terms separately. For (A₁), we have

$$\begin{aligned} (A_1) &= \sum_{s_{t_h+1}} V_{h+1}(\tau_{h+1}) \left((1 - \hat{\theta}(s_{t_h}, a_{t_h}, \delta_{t_h})) \hat{p}_h(s_{t_h+1}|s_{t_h}, a_{t_h}) - (1 - \theta(s_{t_h}, a_{t_h}, \delta_{t_h})) p_h(s_{t_h+1}|s_{t_h}, a_{t_h}) \right) \\ &= \sum_{s_{t_h+1}} V_{h+1}(\tau_{h+1}) \left(\left(1 - \hat{\theta}(s_{t_h}, a_{t_h}, \delta_{t_h}) \right) - (1 - \theta(s_{t_h}, a_{t_h}, \delta_{t_h})) \right) \hat{p}(s_{t_h+1}|s_{t_h}, a_{t_h}) \\ &\quad + \sum_{s_{t_h+1}} V_{h+1}(\tau_{h+1}) (1 - \theta(s_{t_h}, a_{t_h}, \delta_{t_h})) (\hat{p}(s_{t_h+1}|s_{t_h}, a_{t_h}) - p(s_{t_h+1}|s_{t_h}, a_{t_h})) \\ &\stackrel{(i)}{\geq} -H \left| \hat{\theta}(s_{t_h}, a_{t_h}, \delta_{t_h}) - \theta(s_{t_h}, a_{t_h}, \delta_{t_h}) \right| - c_{A,2}H \sqrt{\frac{H\iota}{N_{t_h}(s_{t_h}, a_{t_h})}}, \end{aligned}$$

479 where in (i), the first term is the estimation error of $\hat{\theta}$ using the collected data, the second term follows
 480 from Hoeffding's inequality, and $c_{A,2}$ is an absolute constant. For (A₂), we have

$$(A_2) \geq -H \left| \hat{\theta}(s_{t_h}, a_{t_h}, \delta_{t_h}) - \theta(s_{t_h}, a_{t_h}, \delta_{t_h}) \right|,$$

481 since τ_{h+1} is now uniquely determined. Summing up (A₁) and (A₂), we obtain

$$(A) = (A_1) + (A_2) \geq -2H \left| \hat{\theta}(s_{t_h}, a_{t_h}, \delta_{t_h}) - \theta(s_{t_h}, a_{t_h}, \delta_{t_h}) \right| - c_{A,2}H \sqrt{\frac{H\iota}{N_{t_h}(s_{t_h}, a_{t_h})}}.$$

482 It remains to bound the estimation error of $\widehat{\theta}(s_{t_h}, a_{t_h}, \delta_{t_h})$. Using the Hoeffding's inequality again,
 483 we obtain

$$\left| \widehat{\theta}(s_{t_h}, a_{t_h}, \delta_{t_h}) - \theta(s_{t_h}, a_{t_h}, \delta_{t_h}) \right| \leq c_\theta \sqrt{\frac{H\iota}{N_{t_h}(s_{t_h}, a_{t_h}, \delta_{t_h})}}.$$

484 Taking $c_A = \max\{c_{A,1}, c_{A,2}, c_\theta, 2\}$, we have

$$(A) \geq -c_A H \left(\sqrt{\frac{H\iota}{N_{t_h}(s_{t_h}, a_{t_h}, \delta_{t_h})}} + \sqrt{\frac{H\iota}{N_{t_h}(s_{t_h}, a_{t_h})}} \right).$$

485 With the choice of the bonus function, it can be checked that

$$\widetilde{Q}_{h,\text{aug}}^k(\tau_h, a_h) - \widetilde{Q}_{h,\text{aug}}(\tau_h, a_h) \geq (A) + b_h^k(\tau_h, a_h) \geq 0$$

486 with probability $1 - \gamma$ for any (τ_h, a_h) . \square

487 **Regret analysis** In the sequel, we omit subscripts ‘‘aug’’ and tilde $\widetilde{\cdot}$ for simplicity. Thanks to
 488 Lemma B.1, we consider $(Q_h^k - Q_h^{\pi_k})(\tau_h^k, a_h^k)$ as an upper bound of $(Q_h - Q_h^{\pi_k})(\tau_h^k, a_h^k)$. We
 489 bound $(Q_h^k - Q_h^{\pi_k})(\tau_h^k, a_h^k)$ by

$$\begin{aligned} & (Q_h^k - Q_h^{\pi_k})(\tau_h^k, a_h^k) \\ & \leq \left([\widehat{\mathcal{P}}_h^k V_{h+1}^k - \mathcal{P}_h V_{h+1}^{\pi_k}] \right) (\tau_h^k, a_h^k) + b_h^k(\tau_h^k, a_h^k) \\ & \leq \left([\widehat{\mathcal{P}}_h^k - \mathcal{P}_h] V_{h+1} \right) (\tau_h^k, a_h^k) + \left([\widehat{\mathcal{P}}_h^k - \mathcal{P}_h] [V_{h+1}^k - V_{h+1}] \right) (\tau_h^k, a_h^k) \\ & \quad + \left(\mathcal{P}_h [V_{h+1}^k - V_{h+1}^{\pi_k}] \right) (\tau_h^k, a_h^k) + b_h^k(\tau_h^k, a_h^k) \\ & \leq \underbrace{\left([\widehat{\mathcal{P}}_h^k - \mathcal{P}_h] [V_{h+1}^k - V_{h+1}] \right)}_{(A)} (\tau_h^k, a_h^k) + \left(\mathcal{P}_h [V_{h+1}^k - V_{h+1}^{\pi_k}] \right) (\tau_h^k, a_h^k) + 2b_h^k(\tau_h^k, a_h^k). \end{aligned} \quad (\text{B.1})$$

490 Similar to Lemma B.1, for $h \geq H$, we expand term (A) into

$$\begin{aligned} (A) & = \sum_{\tau_{h+1}} \left(\widehat{p}_h^k(\tau_{h+1} | \tau_h^k, a_h^k) - p_h(\tau_{h+1} | \tau_h^k, a_h^k) \right) [V_{h+1}^k - V_{h+1}](\tau_{h+1}) \\ & = \sum_{s_{t_{h+1}}} [V_{h+1}^k - V_{h+1}](\tau_{h+1}) \left(\widehat{p}_h^k(s_{t_{h+1}} | s_{t_h}^k, a_{t_h}^k) - p_h(s_{t_{h+1}} | s_{t_h}^k, a_{t_h}^k) \right). \end{aligned} \quad (\text{B.2})$$

491 On the other hand, for $h \leq H$, the decomposition of term (A) is more complicated. We have

$$\begin{aligned} (A) & = \sum_{\tau_{h+1}} \left(\widehat{p}_h^k(\tau_{h+1} | \tau_h^k, a_h^k) - p_h(\tau_{h+1} | \tau_h^k, a_h^k) \right) [V_{h+1}^k - V_{h+1}](\tau_{h+1}) \\ & = \underbrace{\sum_{\tau_{h+1}} [V_{h+1}^k - V_{h+1}](\tau_{h+1}) \left(\widehat{p}_h^k(\tau_{h+1} | \tau_h^k, a_h^k) - p_h(\tau_{h+1} | \tau_h^k, a_h^k) \right) \mathbf{1}\{\delta_{t_{h+1}} = 0\} \mathbf{1}\{t_{h+1} = t_h^k + 1\}}_{(A_1)} \\ & \quad + \underbrace{\sum_{\tau_{h+1}} [V_{h+1}^k - V_{h+1}](\tau_{h+1}) \left(\widehat{p}_h^k(\tau_{h+1} | \tau_h^k, a_h^k) - p_h(\tau_{h+1} | \tau_h^k, a_h^k) \right) \mathbf{1}\{\delta_{t_{h+1}} = \delta_{t_h^k} + 1\} \mathbf{1}\{t_{h+1} = t_h^k\}}_{(A_2)}. \end{aligned}$$

492 Term (A₂) can be directly bounded by

$$\begin{aligned} (A_2) & \leq H \left| \widehat{\theta}^k(s_{t_h}^k, a_{t_h}^k, \delta_{t_h}^k) - \theta(s_{t_h}^k, a_{t_h}^k, \delta_{t_h}^k) \right| \\ & \leq c_\theta H \sqrt{\frac{H\iota}{N_{t_h}^k(s_{t_h}^k, a_{t_h}^k, \delta_{t_h}^k)}} \end{aligned}$$

493 with probability $1 - \gamma$. To bound (A_1) , we have

$$\begin{aligned}
(A_1) &= \sum_{s_{t_{h+1}}} [V_{h+1}^k - V_{h+1}] (\tau_{h+1}) \left(\left(1 - \widehat{\theta}^k(s_{t_h}^k, a_{t_h}^k, \delta_{t_h}^k) \right) \widehat{p}_h^k(s_{t_{h+1}} | s_{t_h}^k, a_{t_h}^k) \right. \\
&\quad \left. - (1 - \theta(s_{t_h}^k, a_{t_h}^k, \delta_{t_h}^k)) p_h(s_{t_{h+1}} | s_{t_h}^k, a_{t_h}^k) \right) \\
&= \sum_{s_{t_{h+1}}} [V_{h+1}^k - V_{h+1}] (\tau_{h+1}) \left(\left(1 - \widehat{\theta}^k(s_{t_h}^k, a_{t_h}^k, \delta_{t_h}^k) \right) - (1 - \theta(s_{t_h}^k, a_{t_h}^k, \delta_{t_h}^k)) \right) \widehat{p}_h^k(s_{t_{h+1}} | s_{t_h}^k, a_{t_h}^k) \\
&\quad + \sum_{s_{t_{h+1}}} [V_{h+1}^k - V_{h+1}] (\tau_{h+1}) (1 - \theta(s_{t_h}^k, a_{t_h}^k, \delta_{t_h}^k)) (\widehat{p}_h^k(s_{t_{h+1}} | s_{t_h}^k, a_{t_h}^k) - p_h(s_{t_{h+1}} | s_{t_h}^k, a_{t_h}^k)) \\
&\leq (1 - \theta(s_{t_h}^k, a_{t_h}^k, \delta_{t_h}^k)) \sum_{s_{t_{h+1}}} [V_{h+1}^k - V_{h+1}] (\tau_{h+1}) (\widehat{p}_h^k(s_{t_{h+1}} | s_{t_h}^k, a_{t_h}^k) - p_h(s_{t_{h+1}} | s_{t_h}^k, a_{t_h}^k)) \\
&\quad + H \left| \widehat{\theta}^k(s_{t_h}^k, a_{t_h}^k, \delta_{t_h}^k) - \theta(s_{t_h}^k, a_{t_h}^k, \delta_{t_h}^k) \right| \\
&\leq (1 - \theta(s_{t_h}^k, a_{t_h}^k, \delta_{t_h}^k)) \sum_{s_{t_{h+1}}} [V_{h+1}^k - V_{h+1}] (\tau_{h+1}) (\widehat{p}_h^k(s_{t_{h+1}} | s_{t_h}^k, a_{t_h}^k) - p_h(s_{t_{h+1}} | s_{t_h}^k, a_{t_h}^k)) \\
&\quad + c_\theta H \sqrt{\frac{H\iota}{N_{t_h}^k(s_{t_h}^k, a_{t_h}^k, \delta_{t_h}^k)}}.
\end{aligned}$$

494 Putting (A_1) and (A_2) together, we obtain

$$\begin{aligned}
(A) &\leq (1 - \theta(s_{t_h}^k, a_{t_h}^k, \delta_{t_h}^k)) \sum_{s_{t_{h+1}}} [V_{h+1}^k - V_{h+1}] (\tau_{h+1}) (\widehat{p}_h^k(s_{t_{h+1}} | s_{t_h}^k, a_{t_h}^k) - p_h(s_{t_{h+1}} | s_{t_h}^k, a_{t_h}^k)) \\
&\quad + 2c_\theta H \sqrt{\frac{H\iota}{N_{t_h}^k(s_{t_h}^k, a_{t_h}^k, \delta_{t_h}^k)}}. \tag{B.3}
\end{aligned}$$

495 In both (B.2) and (B.3) for different ranges of h , we apply the Bernstein inequality (Lemma D.1) to
496 derive

$$\begin{aligned}
&\sum_{s_{t_{h+1}}} [V_{h+1}^k - V_{h+1}] (\tau_{h+1}) (\widehat{p}_h^k(s_{t_{h+1}} | s_{t_h}^k, a_{t_h}^k) - p_h(s_{t_{h+1}} | s_{t_h}^k, a_{t_h}^k)) \\
&\leq c \cdot \sum_{s_{t_{h+1}}} [V_{h+1}^k - V_{h+1}] (\tau_{h+1}) \left[\sqrt{\frac{p_h(s_{t_{h+1}} | s_{t_h}^k, a_{t_h}^k) \iota}{N_{t_h}^k(s_{t_h}^k, a_{t_h}^k)}} + \frac{\iota}{N_{t_h}^k(s_{t_h}^k, a_{t_h}^k)} \right] \\
&\stackrel{(i)}{\leq} c \cdot \sum_{s_{t_{h+1}}} [V_{h+1}^k - V_{h+1}] (\tau_{h+1}) \left[\frac{p_h(s_{t_{h+1}} | s_{t_h}^k, a_{t_h}^k)}{2cH} + \frac{(2cH + 1)\iota}{N_{t_h}^k(s_{t_h}^k, a_{t_h}^k)} \right] \\
&\leq c \cdot \left(\frac{SH(2cH + 1)\iota}{N_{t_h}^k(s_{t_h}^k, a_{t_h}^k)} + \frac{1}{2cH} \sum_{s_{t_{h+1}}} [V_{h+1}^k - V_{h+1}] (\tau_{h+1}) p_h(s_{t_{h+1}} | s_{t_h}^k, a_{t_h}^k) \right), \tag{B.4}
\end{aligned}$$

497 where inequality (i) follows from $\sqrt{ab} \leq a + b$. Substituting (B.4) into (B.2), for $h \geq H$, we deduce

$$\begin{aligned}
(A) &\leq \frac{1}{2H} \sum_{s_{t_{h+1}}} [V_{h+1}^k - V_{h+1}] (\tau_{h+1}) p_h(s_{t_{h+1}} | s_{t_h}^k, a_{t_h}^k) + \frac{cSH(2cH + 1)\iota}{N_{t_h}^k(s_{t_h}^k, a_{t_h}^k)} \\
&\stackrel{(i)}{\leq} \frac{1}{2H} (\mathcal{P}_h[V_{h+1}^k - V_{h+1}^{\pi_k}]) (\tau_h^k, a_h^k) + c^\rho \frac{mSH^2\iota}{N_{t_h}^k(s_{t_h}^k, a_{t_h}^k)},
\end{aligned}$$

498 where c^θ is a sufficiently large constant. By the same reasoning, substituting (B.4) into (B.3), for
 499 $h < H$, we have

$$(A) \leq \frac{1}{2H} (1 - \theta(s_{t_h}^k, a_{t_h}^k, \delta_{t_h}^k)) \sum_{s_{t_{h+1}}^k} [V_{h+1}^k - V_{h+1}] (\tau_{h+1}) p_h(s_{t_{h+1}}^k | s_{t_h}^k, a_{t_h}^k) + \frac{cSH(2cH+1)\iota}{N_{t_h}^k(s_{t_h}^k, a_{t_h}^k)}$$

$$+ 2c_\theta H \sqrt{\frac{H\iota}{N_{t_h}^k(s_{t_h}^k, a_{t_h}^k, \delta_{t_h}^k)}}$$

$$\stackrel{(i)}{\leq} \frac{1}{2H} (\mathcal{P}_h[V_{h+1}^k - V_{h+1}^{\pi_k}]) (\tau_h^k, a_h^k) + c^\theta \frac{SH^2\iota}{N_{t_h}^k(s_{t_h}^k, a_{t_h}^k)} + 2c_\theta H \sqrt{\frac{H\iota}{N_{t_h}^k(s_{t_h}^k, a_{t_h}^k, \delta_{t_h}^k)}}.$$

500 We denote $\zeta_h^k = c^\theta \frac{SH^2\iota}{N_{t_h}^k(s_{t_h}^k, a_{t_h}^k)}$. Now we have a unified upper bound on (A) for any $h \in [1, \tilde{H}]$ as

$$(A) \leq \frac{1}{2H} (\mathcal{P}_h[V_{h+1}^k - V_{h+1}^{\pi_k}]) (\tau_h^k, a_h^k) + \zeta_h^k + 2c_\theta H \sqrt{\frac{H\iota}{N_{t_h}^k(s_{t_h}^k, a_{t_h}^k, \delta_{t_h}^k)}}. \quad (\text{B.5})$$

501 Substituting (B.5) back into (B.1), we have

$$(V_h^k - V_h^{\pi_k}) (\tau_h^k) = (Q_h^k - Q_h^{\pi_k}) (\tau_h^k, a_h^k)$$

$$\leq \left(1 + \frac{1}{2H}\right) (\mathcal{P}_h[V_h^k - V_h^{\pi_k}]) (\tau_h^k, a_h^k) + \zeta_h^k + 2b_h^k + 2c_\theta H \sqrt{\frac{H\iota}{N_{t_h}^k(s_{t_h}^k, a_{t_h}^k, \delta_{t_h}^k)}}.$$

502 We further denote $\xi_h^k = (\mathcal{P}_h[V_h^k - V_h^{\pi_k}]) (\tau_h^k, a_h^k) - [V_{h+1}^k - V_{h+1}^{\pi_k}] (\tau_{h+1}^k)$ and rewrite
 503 $(V_h^k - V_h^{\pi_k}) (\tau_h^k)$ as

$$(V_h^k - V_h^{\pi_k}) (\tau_h^k) \leq \left(1 + \frac{1}{2H}\right) ([V_{h+1}^k - V_{h+1}^{\pi_k}] (\tau_{h+1}^k) + \xi_h^k) + \zeta_h^k + 2b_h^k + 2c_\theta H \sqrt{\frac{H\iota}{N_{t_h}^k(s_{t_h}^k, a_{t_h}^k, \delta_{t_h}^k)}}.$$

504 Recall $\tilde{H} = 2H$. Using a recursive summation argument, we deduce

$$(V_1^k - V_1^{\pi_k}) (\tau_1^k) \leq \sum_{h=1}^{\tilde{H}} \left(1 + \frac{1}{2H}\right)^h \left(\xi_h^k + \zeta_h^k + 2b_h^k + 2c_\theta H \sqrt{\frac{H\iota}{N_{t_h}^k(s_{t_h}^k, a_{t_h}^k, \delta_{t_h}^k)}} \right)$$

$$\leq e \sum_{h=1}^{2H} \left(\xi_h^k + \zeta_h^k + 2b_h^k + 2c_\theta H \sqrt{\frac{H\iota}{N_{t_h}^k(s_{t_h}^k, a_{t_h}^k, \delta_{t_h}^k)}} \right).$$

505 As a consequence, the total regret is bounded by

$$\text{Regret}(K) \leq e \sum_{k=1}^K \sum_{h=1}^{2H} \left(\xi_h^k + \zeta_h^k + 2b_h^k + 2c_\theta H \sqrt{\frac{H\iota}{N_{t_h}^k(s_{t_h}^k, a_{t_h}^k, \delta_{t_h}^k)}} \right). \quad (\text{B.6})$$

506 We need to sum over $\zeta_h^k, \xi_h^k, b_h^k$. Consider ζ_h^k first. We have

$$\sum_{k=1}^K \sum_{h=1}^{2H} \zeta_h^k = c^\theta \sum_{k=1}^K \sum_{h=1}^{2H} \frac{SH^2\iota}{N_{t_h}^k(s_{t_h}^k, a_{t_h}^k)}$$

$$\stackrel{(i)}{\leq} c^\theta H \sum_{k=1}^K \sum_{h=1}^H \frac{SH^2\iota}{N_h^k(s_h^k, a_h^k)}$$

$$\stackrel{(ii)}{\leq} c_\zeta H^4 S^2 A \iota^2, \quad (\text{B.7})$$

507 where inequality (i) invokes the fact that t_h only takes value in $\{1, \dots, H\}$ and each $N_{t_h}^k(s_{t_h}^k, a_{t_h}^k)$
 508 is repeated at most H times, and inequality (ii) follows from the pigeon-hole argument in Azar et al.
 509 [2017].

510 Next we bound the summation over ξ_h^k . This is a martingale difference sequence. We apply Azuma-
 511 Hoeffding's inequality (Lemma D.3) with $n = 2H$ and $c_i = 4H$ to obtain

$$\sum_{k=1}^K \sum_{h=1}^{2H} \xi_h^k \leq c_\xi \sqrt{KH^4\iota}. \quad (\text{B.8})$$

512 The additional H dependence above comes from a union bound over $\tilde{\mathcal{S}}_{\text{aug}} \times \mathcal{A}$. Lastly, we tackle the
 513 summation over bonus functions b_h^k . We have

$$\begin{aligned} \sum_{k=1}^K \sum_{h=1}^{2H} b_h^k &= \sum_{k=1}^K \sum_{h=1}^{2H} c_A H \sqrt{\frac{H\iota}{N_{t_h}^k(s_{t_h}, a_{t_h})}} \\ &\leq c_A H \sum_{k=1}^K \sum_{h=1}^H H \sqrt{\frac{H\iota}{N_{t_h}^k(s_{t_h}, a_{t_h})}} \\ &\leq c_b H^{7/2} \sqrt{SAK\iota}. \end{aligned} \quad (\text{B.9})$$

514 Putting (B.7), (B.8) and (B.9) together, we deduce

$$\text{Regret}(K) \leq c \left(H^{7/2} \sqrt{SAK\iota} + H^4 S^2 A \iota^2 + \sqrt{H^4 K \iota} \right) + 2ec_\theta H \sum_{k=1}^K \sum_{h=1}^{2H} \sqrt{\frac{H\iota}{N_{t_h}^k(s_{t_h}^k, a_{t_h}^k, \delta_{t_h}^k)}}$$

515 for some constant c . To this end, the only remaining task is to find $\sum_{k=1}^K \sum_{h=1}^{2H} \sqrt{\frac{1}{N_{t_h}^k(s_{t_h}^k, a_{t_h}^k, \delta_{t_h}^k)}}$,

516 which undergoes a similar argument as the bonus summation. We have

$$\begin{aligned} \sum_{k=1}^K \sum_{h=1}^{2H} \sqrt{\frac{1}{N_{t_h}^k(s_{t_h}^k, a_{t_h}^k, \delta_{t_h}^k)}} &\leq H \sum_{k=1}^K \sum_{h=1}^H \sqrt{\frac{1}{N_h^k(s_h^k, a_h^k, \delta_h^k)}} \\ &= H \sum_{(h,s,a,\delta)} N_h^K(s,a,\delta) \sum_{i=1}^H \sqrt{\frac{1}{i}} \\ &\stackrel{(i)}{\leq} 2H \sum_{\delta} \sum_{(h,s,a)} \sqrt{N_h^K(s,a,\delta)} \\ &\stackrel{(ii)}{\leq} 2H \sum_{\delta} \sqrt{SAKH} \\ &\stackrel{(iii)}{\leq} 2H^2 \sqrt{SAKH}, \end{aligned} \quad (\text{B.10})$$

517 where inequality (i) invokes $\sum_{i=1}^n 1/\sqrt{i} \leq 2\sqrt{n}$, inequality (ii) follows from Cauchy-Schwarz, and
 518 inequality (iii) uses the fact that δ is bounded by H . Plugging (B.10) into the regret bound, we obtain
 519 the desired result

$$\text{Regret}(K) \leq c \left(H^4 \sqrt{SAK\iota} + H^4 S^2 A \iota^2 + \sqrt{H^4 K \iota} \right)$$

520 with probability $1 - \gamma$. Absorbing $\sqrt{H^4 K \iota}$ into $H^4 \sqrt{SAK\iota}$ yields the bound in Theorem 4.1. \square

521 B.2 Statement and proof of Proposition B.2

522 **Proposition B.2.** In the setup of Section 4.2, we have

$$\begin{aligned} \text{gap}(s_1) &\leq \sum_{h=1}^H \left[\underbrace{\int_{\tau} \left(\mathbb{E}_{s \sim b_h(\cdot|\tau)} [\max_a r_h(s,a)] - \max_a \mathbb{E}_{s \sim b_h(\cdot|\tau)} [r_h(s,a)] \right) \left(\rho_h^{\pi_{\text{delay}}} \wedge \rho_h^{\pi_{\text{nodelay}}} \right) (\tau) d\tau}_{\mathcal{E}_1} \right. \\ &\quad \left. + 2 \underbrace{\|\rho_h^{\pi_{\text{nodelay}}} - \rho_h^{\pi_{\text{delay}}}\|_{\text{TV}}}_{\mathcal{E}_2} \right]. \end{aligned}$$

523 where $\rho_h^{\pi_{\text{nodelay}}}$ and $\rho_h^{\pi_{\text{delay}}}$ are visitation measures induced by π_{nodelay} and π_{delay} , respectively.

524 Term \mathcal{E}_1 is strictly larger than zero due to the convexity of the max operation. Term \mathcal{E}_2 accounts
 525 for the difference in the visitation measure. When the original MDP has deterministic transitions,
 526 we can check that \mathcal{E}_1 is zero, since the expectation over s is concentrated on a singleton that
 527 can be inferred from history. Hence, the visitation measures are also identical, which implies
 528 $V_{1,\text{nodelay}}(s_1) - V_{1,\text{delay}}(s_1) = 0$. On the contrary, when $b_h(\cdot|\tau)$ is evenly spread, meaning that
 529 the entropy of b_h is high, we potentially suffer from a large performance drop, in that, inferring the
 530 current state is difficult.

531 *Proof of Proposition B.2.* Let τ_1, \dots, τ_H denote the states observed in the delayed environment.
 532 Since π_{nodelay} is greedy and Markov, we obtain

$$\begin{aligned} V_{1,\text{nodelay}}(s_1) &= \mathbb{E}^{\pi_{\text{nodelay}}} \left[\sum_{h=1}^{H-1} r_h(s_h, a_h) \mid s_1 \right] + \mathbb{E}^{\pi_{\text{nodelay}}} [E[r_H(s_H, a_H) \mid \tau_H] \mid s_1] \\ &= \mathbb{E}^{\pi_{\text{nodelay}}} \left[\sum_{h=1}^{H-1} r_h(s_h, a_h) \mid s_1 \right] + \mathbb{E}^{\pi_{\text{nodelay}}} \left[\sum_s \mathfrak{b}_H(s \mid \tau_H) \max_a r_H(s, a) \mid s_1 \right]. \end{aligned}$$

533 Recursively applying the above argument, we deduce

$$V_{1,\text{nodelay}}(s_1) = \mathbb{E}^{\pi_{\text{nodelay}}} \left[\sum_{h=1}^H \sum_s \mathfrak{b}_h(s \mid \tau_h) \max_a r_h(s, a) \mid s_1 \right].$$

534 We also rewrite $V_{1,\text{delay}}(s_1)$ as

$$\begin{aligned} V_{1,\text{delay}}(s_1) &= \mathbb{E}^{\pi_{\text{delay}}} \left[\sum_{h=1}^{H-1} r_h(s_h, a_h) \mid s_1 \right] + \mathbb{E}^{\pi_{\text{delay}}} [E[r_H(s_H, a_H) \mid \tau_H] \mid s_1] \\ &= \mathbb{E}^{\pi_{\text{delay}}} \left[\sum_{h=1}^{H-1} r_h(s_h, a_h) \mid s_1 \right] + \mathbb{E}^{\pi_{\text{delay}}} \left[\max_a \sum_s \mathfrak{b}_H(s \mid \tau_H) r_H(s, a) \mid s_1 \right] \\ &= \dots \\ &= \mathbb{E}^{\pi_{\text{delay}}} \left[\sum_{h=1}^H \max_a \sum_s \mathfrak{b}_h(s \mid \tau_h) r_h(s, a) \mid s_1 \right]. \end{aligned}$$

535 Then we write the difference between $V_{1,\text{nodelay}}(s_1)$ and $V_{1,\text{delay}}(s_1)$ as

$$\begin{aligned} &V_{1,\text{nodelay}}(s_1) - V_{1,\text{delay}}(s_1) \\ &= \sum_{h=1}^H \left(\int_{\tau} \sum_s \max_a \mathfrak{b}_h(s \mid \tau) r_h(s, a) \rho_h^{\pi_{\text{nodelay}}}(\tau) d\tau - \int_{\tau} \max_a \sum_s \mathfrak{b}_h(s \mid \tau) r_h(s, a) \rho_h^{\pi_{\text{delay}}}(\tau) d\tau \right) \\ &= \sum_{h=1}^H \left(\int_{\tau} \sum_s \max_a \mathfrak{b}_h(s \mid \tau) r_h(s, a) \rho_h^{\pi_{\text{nodelay}}}(\tau) d\tau - \int_{\tau} \max_a \sum_s \mathfrak{b}_h(s \mid \tau) r_h(s, a) \rho_h^{\pi_{\text{nodelay}}}(\tau) d\tau \right. \\ &\quad \left. + \int_{\tau} \max_a \sum_s \mathfrak{b}_h(s \mid \tau) r_h(s, a) \rho_h^{\pi_{\text{nodelay}}}(\tau) d\tau - \int_{\tau} \max_a \sum_s \mathfrak{b}_h(s \mid \tau) r_h(s, a) \rho_h^{\pi_{\text{delay}}}(\tau) d\tau \right) \\ &\leq \sum_{h=1}^H \left[\int_{\tau} \left(\mathbb{E}_s \mathfrak{b}_h(j\tau) [\max_a r_h(s, a)] - \max_a \mathbb{E}_s \mathfrak{b}_h(j\tau) [r_h(s, a)] \right) \rho_h^{\pi_{\text{nodelay}}}(\tau) d\tau + 2 \|\rho_h^{\pi_{\text{nodelay}}} - \rho_h^{\pi_{\text{delay}}}\|_{\text{TV}} \right]. \end{aligned}$$

536 We also have

$$\begin{aligned} &V_{1,\text{nodelay}}(s_1) - V_{1,\text{delay}}(s_1) \\ &= \sum_{h=1}^H \left(\int_{\tau} \sum_s \max_a \mathfrak{b}_h(s \mid \tau) r_h(s, a) \rho_h^{\pi_{\text{nodelay}}}(\tau) d\tau - \int_{\tau} \sum_s \max_a \mathfrak{b}_h(s \mid \tau) r_h(s, a) \rho_h^{\pi_{\text{delay}}}(\tau) d\tau \right. \\ &\quad \left. + \int_{\tau} \sum_s \max_a \mathfrak{b}_h(s \mid \tau) r_h(s, a) \rho_h^{\pi_{\text{delay}}}(\tau) d\tau - \int_{\tau} \max_a \sum_s \mathfrak{b}_h(s \mid \tau) r_h(s, a) \rho_h^{\pi_{\text{delay}}}(\tau) d\tau \right) \\ &\leq \sum_{h=1}^H \left[\int_{\tau} \left(\mathbb{E}_s \mathfrak{b}_h(j\tau) [\max_a r_h(s, a)] - \max_a \mathbb{E}_s \mathfrak{b}_h(j\tau) [r_h(s, a)] \right) \rho_h^{\pi_{\text{delay}}}(\tau) d\tau + 2 \|\rho_h^{\pi_{\text{nodelay}}} - \rho_h^{\pi_{\text{delay}}}\|_{\text{TV}} \right]. \end{aligned}$$

537 Combining the above two inequalities, we obtain

$$\begin{aligned} & V_{1,\text{nodelay}}(s_1) - V_{1,\text{delay}}(s_1) \\ & \leq \sum_{h=1}^H \left[\int_{\tau} \left(\mathbb{E}_{s \sim b_h(\cdot|\tau)} [\max_a r_h(s, a)] - \max_a \mathbb{E}_{s \sim b_h(\cdot|\tau)} [r_h(s, a)] \right) \left(\rho_h^{\pi_{\text{delay}}} \wedge \rho_h^{\pi_{\text{nodelay}}} \right) (\tau) d\tau \right. \\ & \quad \left. + 2 \|\rho_h^{\pi_{\text{nodelay}}} - \rho_h^{\pi_{\text{delay}}}\|_{\text{TV}} \right]. \end{aligned}$$

538 The proof is complete. \square

539 B.3 Proof of Proposition 4.2

540 *Proof of Proposition 4.2.* We construct an MDP instance $(\mathcal{S}, \mathcal{A}, H, R, P)$ for $H > d$ as follows. Let
541 $\mathcal{S} = \{1, 2\}$ and $\mathcal{A} = \{a_1, a_2\}$. For the reward function, we have

$$r_h(s, a) = \begin{cases} 1 & \text{if } a = a_s \text{ and } h = d + 1 \\ 0 & \text{otherwise} \end{cases}.$$

542 The reward is nonzero only at time $d + 1$. The transition probabilities are defined as

$$p_h(s^\ell | s, a) = \begin{cases} \frac{1}{2} & \text{if } h = d + 1 \\ 1 & \text{if } h \neq d + 1 \text{ and } s^\ell = s \\ 0 & \text{otherwise} \end{cases}.$$

543 The transition probability at step $d + 1$ says that s^ℓ is uniform regardless of the previous state and
544 action. Suppose a uniform initial distribution on s_1 . We first show that if the constant delay equals d ,
545 then there exists a policy π^d achieving maximal value. Indeed, the policy is chosen as

$$\pi_h^d(\cdot | \{s_{h-d}, \mathbf{a}_{h-d:h-1}\}) = \begin{cases} a_{s_{h-d}} & \text{if } h = d + 1 \\ \text{Uniform}(\mathcal{A}) & \text{if } h \neq d + 1. \end{cases}$$

546 It is straightforward to check that π^d is optimal, since at step $d + 1$, s_1 is revealed and the policy
547 takes the optimal action a_{s_1} to obtain reward 1.

548 On the other hand, if the constant delay equals $d + 1$, then any policy suffers from a constant
549 performance degradation. To see this, in a single trajectory, since the starting state is only revealed at
550 time $d + 2$, the policy at time $d + 1$ cannot exploit the information of the initial state. Therefore, any
551 policy coincides with the best action with probability $\frac{1}{2}$. For K episodes, with probability $1 - \gamma$, the
552 total reward of any policy $\pi \in \Pi_e$ is bounded by

$$\sum_{k=1}^K V_1^\pi(s_1^k) \leq \frac{1}{2}K + \sqrt{\frac{K}{2} \log \frac{1}{\gamma}},$$

553 due to Hoeffding's inequality. As a result, the performance drop is at least

$$\text{gap}(K) \geq \frac{1}{2} - \sqrt{\frac{1}{2K} \log \frac{1}{\gamma}}.$$

554 \square

555 **C Omitted proofs in Section 5**

556 **C.1 Proof of Proposition 5.1**

557 *Proof of Proposition 5.1.* We have by standard performance difference arguments that

$$\begin{aligned}
\sum_{k=1}^K \max_{\pi \in \Pi_e} V_{\theta^*}^{\pi^k}(s_1^k) - V_{\theta^*}^{\pi^k}(s_1^k) &\stackrel{(i)}{\leq} \sum_{k=1}^K V_{\theta^*}^{\pi^k}(s_1^k) - V_{\theta^*}^{\pi^k}(s_1^k) \\
&\stackrel{(ii)}{=} \sum_{k=1}^K \sum_{h=1}^H \mathbb{E}_{\theta^*}^{\pi^k} \left[\left\langle (P_h^{\theta^k} - P_h^{\theta^*})(\cdot | s_h, a_h), V_{\theta^*}^{\pi^k, h+1}(\cdot) \right\rangle \right] \\
&\leq \sum_{h=1}^H \sum_{k=1}^K \mathbb{E}_{\theta^*}^{\pi^k} \left[c \sqrt{\frac{H^2 S \iota}{N_h^k(s_h, a_h)}} \wedge H \right] \\
&\stackrel{(iii)}{\leq} \sum_{h=1}^H \sum_{k=1}^K c^\rho \sqrt{\frac{H^2 S \iota}{N_h^k(s_h^k, a_h^k)}} + H \sqrt{H^2 K \iota} \\
&\stackrel{(iv)}{\leq} c^\rho \left(\left[\frac{\log \frac{HK}{\gamma}}{-\log(1 - \lambda_0^2)} \right] \sqrt{H^2 S \iota \cdot S A H K} + \sqrt{H^4 K \iota} \right) \\
&\leq c^\rho \left(\left[\frac{1}{-\log(1 - \lambda_0^2)} \right] \sqrt{H^3 S^2 A K \iota^3} + \sqrt{H^4 K \iota} \right),
\end{aligned}$$

558 where inequality (i) follows from valid optimism, equality (ii) recursively expand the value function
559 and $\langle \cdot, \cdot \rangle$ denotes the inner product, inequality (iii) invokes Azuma-Hoeffding's inequality, and
560 inequality (iv) invokes Lemma C.2. \square

561 **C.2 Algorithm and proof of Theorem 5.2**

Algorithm 4 Policy learning for MDPs with missing observations

- 1: **Input:** Horizon H .
- 2: **Init:** $V_{H+1}(\tau) = 0$ and $Q_H(\tau, a) = H$ for any τ, a , data set $\mathcal{D}^0 = \emptyset$, initial policy π^0 .
- 3: **for** episode $k = 1, \dots, K$ **do**
- 4: Execute policy π^k .
- 5: After the episode ends, collect data $\mathcal{D}^k = \mathcal{D}^{k-1} \cup \{(s_h^k, a_h^k, r_h^k)\}_{h=1}^H$.
- 6: On data set \mathcal{D}^k , compute counting numbers

$$N_h^k(\tau_h, a_h) = \sum_{j=1}^k \mathbb{1}\{\tau_h^j = \tau_h, a_h^j = a_h, s_{h+1}^j \neq \emptyset\} \quad \text{and} \quad N_{h,\lambda}^k = \sum_{j=1}^k \mathbb{1}\{s_h^j = \emptyset\}.$$

- 7: Estimate transition probabilities and delay distributions via

$$\hat{p}_h^k(s_{h+1} | \tau_h, a_h) = \frac{N_h^k(\tau_h, a_h, s_{h+1})}{N_h^k(\tau_h, a_h)} \quad \text{and} \quad \hat{\lambda}_h^k = N_{h,\lambda}^k / k.$$

- 8: Set bonus function as

$$b_h^k(\tau_h, a_h) = cH \left(\sqrt{\frac{H \iota}{N_h^k(\tau_h, a_h)}} + \sqrt{\frac{\iota}{k}} \right)$$

for $\iota = \log \frac{S A K H}{\gamma}$ and c sufficiently large.

- 9: Run optimistic value iteration in MDP_{aug} for H steps and obtain $\pi^k \in \Pi_e$.
 - 10: **end for**
 - 11: **Return:** Learned policy π^k for $k = 1, \dots, K$.
-

562 We remark that similar to delayed MDPs, in Line 9 the planning is on MDP_{aug} and the obtained
563 policy is executable given any $\tau \in \mathcal{S}_{\text{aug}}$ when state observation is missed. Therefore, the planning

564 complexity is SA^H . Different from Algorithm 2, the bonus function here depends on multi-step
 565 transitions, in that missing observations are permanently lost.

566 *Proof of Theorem 5.2.* The proof utilizes similar steps as Theorem 4.1, with an extra care on the
 567 summation of bonus functions.

568 **Valid optimism** We verify the choice of bonus functions leads to a valid optimism.

569 **Lemma C.1.** Given any failure probability $\gamma < 1$, we set bonus functions as

$$b_h^k(\tau_h, a_h) = cH \left(\sqrt{\frac{H\iota}{N_h^k(\tau_h, a_h)}} + \sqrt{\frac{\iota}{k}} \right) \quad \text{with} \quad \iota = \log \left(\frac{SAHK}{\gamma} \right).$$

570 Then with probability $1 - \gamma$, it holds

$$Q_{h,\text{aug}}^k(\tau_h, a_h) \geq Q_{h,\text{aug}}(\tau_h, a_h), \quad V_{h,\text{aug}}^k(\tau_h) \geq V_{h,\text{aug}}(\tau_h) \quad \text{for any} \quad (k, h, \tau_h, a_h).$$

571 *Proof.* In the proof, we omit subscript ‘‘aug’’ for simplicity. We use backward induction on time h
 572 again. The base case of $H + 1$ holds immediately due to the initial value of $V_{H+1,\text{aug}}$. Suppose at time
 573 $h + 1$, the assertion holds. Then for time h , if $Q_{h,\text{aug}} = H$, the assertion holds trivially. Otherwise,
 574 we have

$$\begin{aligned} & Q_h(\tau_h, a_h) - Q_h(\tau_h, a_h) \\ &= \widehat{r}_h(\tau_h, a_h) + [\widehat{\mathcal{P}}_h V_{h+1}](\tau_h, a_h) - r_h(\tau_h, a_h) - [\mathcal{P}_h V_{h+1}](\tau_h, a_h) + b_h^k(\tau_h, a_h) \\ &\geq \underbrace{([\widehat{\mathcal{P}}_h - \mathcal{P}_h] V_{h+1})}_{(A)}(\tau_h, a_h) + \underbrace{\widehat{r}_h(\tau_h, a_h) - r_h(\tau_h, a_h)}_{(B)} + b_h^k(\tau_h, a_h). \end{aligned}$$

575 We lower bound (A) and (B) separately. For term (A), we have

$$\begin{aligned} (A) &= \sum_{\tau_{h+1}} V_{h+1}(\tau_{h+1}) (\widehat{p}_h(\tau_{h+1}|\tau_h, a_h) - p_h(\tau_{h+1}|\tau_h, a_h)) \\ &= \sum_{\tau_{h+1}} V_{h+1}(\tau_{h+1}) (\widehat{p}_h(\tau_{h+1}|\tau_h, a_h) - p_h(\tau_{h+1}|\tau_h, a_h)) \mathbf{1}\{t_{h+1} = h + 1\} \\ &\quad + \sum_{\tau_{h+1}} V_{h+1}(\tau_{h+1}) (\widehat{p}_h(\tau_{h+1}|\tau_h, a_h) - p_h(\tau_{h+1}|\tau_h, a_h)) \mathbf{1}\{t_{h+1} = t_h\} \\ &= \underbrace{\sum_{s_{h+1}} V_{h+1}(\tau_{h+1}) \left((1 - \widehat{\lambda}_h) \widehat{p}_h(s_{h+1}|s_{t_h}, \mathbf{a}_{t_h:h}) - (1 - \lambda_h) p_h(s_{h+1}|s_{t_h}, \mathbf{a}_{t_h:h}) \right)}_{(A_1)} \\ &\quad + \underbrace{V_{h+1}(\{s_{t_h}, \mathbf{a}_{t_h:h}\})}_{(A_2)} (\widehat{\lambda}_h - \lambda_h). \end{aligned}$$

576 In (A₁), τ_{h+1} is $\{s_{h+1}\}$. We bound (A₁) as

$$\begin{aligned} (A_1) &= \sum_{s_{h+1}} V_{h+1}(\tau_{h+1}) \left((1 - \widehat{\lambda}_h) \widehat{p}_h(s_{h+1}|s_{t_h}, \mathbf{a}_{t_h:h}) - (1 - \lambda_h) \widehat{p}_h(s_{h+1}|s_{t_h}, \mathbf{a}_{t_h:h}) \right. \\ &\quad \left. + (1 - \lambda_h) \widehat{p}_h(s_{h+1}|s_{t_h}, \mathbf{a}_{t_h:h}) - (1 - \lambda_h) p_h(s_{h+1}|s_{t_h}, \mathbf{a}_{t_h:h}) \right) \\ &= \sum_{s_{h+1}} V_{h+1}(\tau_{h+1}) (1 - \lambda_h) (\widehat{p}_h(s_{h+1}|s_{t_h}, \mathbf{a}_{t_h:h}) - p_h(s_{h+1}|s_{t_h}, \mathbf{a}_{t_h:h})) \\ &\quad + \sum_{s_{h+1}} V_{h+1}(\tau_{h+1}) (\lambda_h - \widehat{\lambda}_h) \widehat{p}_h(s_{h+1}|s_{t_h}, \mathbf{a}_{t_h:h}) \\ &\stackrel{(i)}{\geq} -c_A H \sqrt{\frac{H\iota}{N_h(\tau_h, a_h)}} - H |\widehat{\lambda}_h - \lambda_h|, \end{aligned}$$

577 where inequality (i) invokes Hoeffding's inequality and holds with probability $1 - \gamma$ for any τ_h, a_h
 578 and some constant c_A . Term (A_2) is immediately bounded by

$$(A_2) \geq -H \left| \widehat{\lambda}_h - \lambda_h \right|.$$

579 Putting (A_1) and (A_2) together, we derive

$$(A) \geq -c_A H \sqrt{\frac{H\iota}{N_h(\tau_h, a_h)}} - 2H \left| \widehat{\lambda}_h - \lambda_h \right|$$

580 with high probability. For term (B) , we have

$$(B) = \sum_{s_h} r(s_h, a_h) \left(\widehat{b}_h(s_h | \tau_h) - b_h(s_h | \tau_h) \right) \geq -c_B \sqrt{\frac{H\iota}{N_h(\tau_h, a_h)}}.$$

581 Taking $c = c_A + c_B$ and summing up (A) and (B) , we have

$$Q_h(\tau_h, a_h) - Q_h(\tau_h, a_h) \geq -cH \sqrt{\frac{H\iota}{N_h(\tau_h, a_h)}} - 2H \left| \widehat{\lambda}_h - \lambda_h \right| + b_h^k(\tau_h, a_h).$$

582 We estimate λ_h by its empirical average. In episode $k \geq 1$, we have access to k i.i.d. realizations of
 583 Bernoulli random variable with rate λ_h (observable or not). Therefore, by Hoeffding's inequality, we
 584 have

$$\left| \widehat{\lambda}_h^k - \lambda_h \right| \leq 2\sqrt{\frac{\log \frac{HK}{\gamma}}{k}} \leq 2\sqrt{\frac{\iota}{k}}.$$

585 Substituting into $Q_h^k(\tau_h, a_h) - Q_h(\tau_h, a_h)$ and reloading constant c give rise to

$$Q_h^k(\tau_h, a_h) - Q_h(\tau_h, a_h) \geq -cH \left(\sqrt{\frac{H\iota}{N_h^k(\tau_h, a_h)}} + \sqrt{\frac{\iota}{k}} \right) + b_h^k(\tau_h, a_h) \geq 0.$$

586 The proof is complete. \square

587 **Regret analysis** We omit subscript "aug" to ease the presentation. The same derivation in the proof
 588 of Theorem 4.1 gives rise to

$$\begin{aligned} (Q_h - Q_h^{\pi_k})(\tau_h^k, a_h^k) &\leq (Q_h^k - Q_h^{\pi_k})(\tau_h^k, a_h^k) \\ &\leq \underbrace{\left([\widehat{\mathcal{P}}_h^k - \mathcal{P}_h][V_{h+1}^k - V_{h+1}] \right)}_{(A)} (\tau_h^k, a_h^k) + (\mathcal{P}_h[V_{h+1}^k - V_{h+1}^{\pi_k}]) (\tau_h^k, a_h^k) + 2b_h^k(\tau_h^k, a_h^k). \end{aligned} \quad (\text{C.1})$$

589 Lemma C.1 shows that (A) can be written as

$$\begin{aligned} (A) &= \sum_{s_{h+1}} [V_{h+1}^k - V_{h+1}](\tau_{h+1})(1 - \lambda_h) \left(\widehat{p}_h^k(s_{h+1} | s_{t_h}^k, \mathbf{a}_{t_h:h}^k) - p_h(s_{h+1} | s_{t_h}^k, \mathbf{a}_{t_h:h}^k) \right) \\ &\quad + \sum_{s_{h+1}} [V_{h+1}^k - V_{h+1}](\tau_{h+1})(\lambda_h - \widehat{\lambda}_h^k) \widehat{p}_h^k(s_{h+1} | s_{t_h}^k, \mathbf{a}_{t_h:h}^k) \\ &\leq \sum_{s_{h+1}} [V_{h+1}^k - V_{h+1}](\tau_{h+1})(1 - \lambda_h) \left(\widehat{p}_h^k(s_{h+1} | s_{t_h}^k, \mathbf{a}_{t_h:h}^k) - p_h(s_{h+1} | s_{t_h}^k, \mathbf{a}_{t_h:h}^k) \right) + H \left| \widehat{\lambda}_h^k - \lambda_h \right| \\ &\leq (1 - \lambda_h) \sum_{s_{h+1}} [V_{h+1}^k - V_{h+1}](\tau_{h+1}) \left(\widehat{p}_h^k(s_{h+1} | s_{t_h}^k, \mathbf{a}_{t_h:h}^k) - p_h(s_{h+1} | s_{t_h}^k, \mathbf{a}_{t_h:h}^k) \right) + 2H \sqrt{\frac{\iota}{k}}. \end{aligned}$$

590 Following the derivation in (B.4), (B.5) and (B.6), we have

$$\begin{aligned} \text{Regret}(K) &\leq e \sum_{k=1}^K \sum_{h=1}^H \left(\xi_h^k + \zeta_h^k + 2b_h^k + 2H \sqrt{\frac{\iota}{k}} \right) \\ &\leq e \sum_{k=1}^K \sum_{h=1}^H \left(\xi_h^k + \zeta_h^k + 2b_h^k \right) + 2\sqrt{H^4 K \iota}. \end{aligned}$$

591 where $\xi_h^k = (\mathcal{P}_h[V_{h+1}^k - V_{h+1}^{\pi_k}])(\tau_h^k, a_h^k) - [V_{h+1}^k - V_{h+1}^{\pi_k}](\tau_{h+1}^k)$ is the martingale difference and

592 $\zeta_h^k = c^\theta \frac{SH^2 \iota}{N_h^k(\tau_h^k, a_h^k)}.$

593 **Counting number summation** The summation over ξ_h^k is standard. Using the Azuma-Hoeffding's
 594 inequality, we have

$$\sum_{k=1}^K \sum_{h=1}^H \xi_h^k \leq c_\xi \sqrt{KH^4\iota}.$$

595 It remains to find the summations involving $N_h^k(\tau_h^k, a_h^k)$. First, we show that the event $\mathcal{E}_m =$
 596 $\{h - t_h - 1 \leq m\}$, i.e., the maximal consecutive delay is upper bounded by $m > 0$, holds with high
 597 probability. We have

$$\mathbb{P}(\mathcal{E}_m) \leq (1 - H(1 - \lambda_0)^{m+1})^K,$$

598 since λ_0 is a uniform lower bound of λ_h . Next, we provide an upper bound on $N_h^K(\tau_h, a_h)$. For
 599 a given tuple (h, τ_h, a_h, t_h) , the consecutive missing length is $h - t_h - 1$. Such a missing pattern
 600 appears with probability at most $(1 - \lambda_0)^{h - t_h - 1}$. As a consequence, denote $C_{h - t_h - 1}^K$ as the number
 601 of $h - t_h - 1$ consecutive missings in K episodes. With probability $1 - \gamma$, we have

$$C_{h - t_h - 1}^K \leq K(1 - \lambda_0)^{h - t_h - 1} + \sqrt{K(1 - \lambda_0)^{h - t_h - 1}H\iota + \iota}.$$

602 by Bernstein's inequality in Lemma D.1. Furthermore, at a fixed time h , we use Lemma C.3 to
 603 bound the gap between two consecutive appearances of the same missing pattern. We instantiate
 604 Lemma C.3 with $\theta = (1 - \lambda_0)^{h - t_h - 1}$ and obtain that the gap is bounded by $\left\lceil \frac{\iota}{\log(1 - (1 - \lambda_0)^{h - t_h - 1})} \right\rceil$
 605 with probability $1 - \gamma$. Within the gap, the number of consecutive delays of length larger than
 606 $h - t_h - 1$ is bounded by

$$\begin{aligned} C_{h - t_h - 1} &\stackrel{(i)}{\leq} \left\lceil \frac{\iota}{-\log(1 - (1 - \lambda_0)^{h - t_h - 1})} \right\rceil (1 - \lambda_0)^{h - t_h} \\ &\quad + \sqrt{\left\lceil \frac{\iota}{-\log(1 - (1 - \lambda_0)^{h - t_h - 1})} \right\rceil (1 - \lambda_0)^{h - t_h} H\iota + \iota} \\ &\stackrel{(ii)}{\leq} \sqrt{2(1 - \lambda_0)H\iota} + 2(1 - \lambda_0) + \iota, \end{aligned}$$

607 where inequality (i) follows from Bernstein's inequality again and inequality (ii) invokes the fact
 608 $x + \log(1 - x) \leq 0$ for $x \in [0, 1)$ and bounds $\lceil x \rceil$ by $x + 1$. Now we can bound the summation of
 609 the counting numbers. Conditioned on the event \mathcal{E}_m , we have

$$\begin{aligned} \sum_{k=1}^K \sum_{h=1}^H \sqrt{\frac{1}{N_h^k(\tau_h^k, a_h^k)}} &\stackrel{(i)}{\leq} \sum_{(h, \tau, a, t_h)} C_{h - t_h - 1} \sum_{i=1}^{N_h^K(\tau, a)} \sqrt{\frac{1}{i}} \\ &\leq 2 \left(\sqrt{2(1 - \lambda_0)H\iota} + 2(1 - \lambda_0) + \iota \right) \sum_{(h, \tau, a, t_h)} \sqrt{N_h^K(\tau, a)} \\ &\stackrel{(ii)}{\leq} 2 \left(\sqrt{2(1 - \lambda_0)H\iota} + 2(1 - \lambda_0) + \iota \right) \sum_{h, t_h} \sqrt{SA^{h - t_h} C_{h - t_h - 1}^K} \\ &\leq 2 \left(\sqrt{2(1 - \lambda_0)H\iota} + 2(1 - \lambda_0) + \iota \right) \\ &\quad \cdot \sum_{h, t_h} \sqrt{SA \left(K((1 - \lambda_0)A)^{h - t_h - 1} + \sqrt{K(A^2(1 - \lambda_0))^{h - t_h - 1}H\iota} + A^{h - t_h - 1}\iota \right)} \\ &\stackrel{(iii)}{\leq} 2 \left(\sqrt{2(1 - \lambda_0)H\iota} + 2(1 - \lambda_0) + \iota \right) \sum_{h, t_h} \sqrt{SA \left(K + \sqrt{KA^m H\iota} + A^m \iota \right)} \\ &\leq 2 \left(\sqrt{2(1 - \lambda_0)H\iota} + 2(1 - \lambda_0) + \iota \right) H^2 \sqrt{SA \left(K + \sqrt{KA^m H\iota} + A^m \iota \right)} \\ &\leq 2\sqrt{H^5 SA\iota^2 \left(K + \sqrt{KA^m H\iota} + A^m \iota \right)}, \end{aligned}$$

610 where inequality (i) follows since N_h^k is repeated at most C_{h-t_h-1} times before getting an up-
611 date and inequality (ii) follows from Cauchy-Schwarz inequality, and inequality (iii) invokes the
612 assumption of $\lambda A \leq 1$. Moreover, conditioned on the event \mathcal{E}_m , we also have

$$\begin{aligned} \sum_{k=1}^K \sum_{h=1}^H \frac{1}{N_h^k(\tau_h^k, a_h^k)} &\leq \sum_{(h,\tau,a,t_h)} C_{h-t_h-1} \sum_{i=1}^{N_h^K(\tau,a)} \frac{1}{i} \\ &\leq \left(\sqrt{2(1-\lambda_0)H\iota} + 2(1-\lambda_0) + \iota \right) \sum_{(h,\tau,a,t_h)} \log N_h^K(\tau, a) \\ &\leq \iota H^{5/2} S A^{m+1} \log K. \end{aligned}$$

613 **Putting together** On event \mathcal{E}_m , the regret is bounded by

$$\begin{aligned} \text{Regret}(K) &\stackrel{(i)}{\leq} c \left(\sqrt{H^4 K \iota} + \sum_{k=1}^K \sum_{h=1}^H \left[\frac{S H^2 \iota}{N_h^k(\tau_h^k, a_h^k)} + H \sqrt{\frac{H \iota}{N_h^k(\tau_h^k, a_h^k)}} \right] \right) \\ &\leq c \left(H^4 \sqrt{S A \iota^3 K} \left(1 + \sqrt{\frac{A^m H \iota}{K}} + \frac{A^m \iota}{K} \right) + S^2 A^m \sqrt{H^9 \iota^6} + \sqrt{H^4 K \iota} \right), \end{aligned}$$

614 where c is a sufficiently large constant and we substitute the bonus functions into inequality (i).

615 On the complement of \mathcal{E}_m , the regret is bounded by $H(1 - \mathbb{P}(\mathcal{E}_m)) \leq H^2 K(1 - \lambda_0)^{m+1}$. We choose

616 $m = \frac{1}{2} \left\lceil \frac{\log K}{-\log(1 - \lambda_0)} \right\rceil$ such that $H(1 - \mathbb{P}(\mathcal{E}_m)) \leq H^2 K(1 - \lambda_0)^{m+1} \leq H^2 \sqrt{K}$. We can now check

617 that $A^{m+1} = \exp\left(-\frac{\log A}{-\log(1 - \lambda_0)} \log \sqrt{K}\right) \leq K^{\frac{1}{2(1+v)}}$. Therefore, combining the regret on event \mathcal{E}_m

618 and the complement event \mathcal{E}_m^c leads to

$$\text{Regret}(K) \leq c \left(H^4 \sqrt{S A K \iota^3} + S^2 \sqrt{H^9 K^{\frac{1}{1+v}} \iota^6} \right).$$

619 The proof is complete. \square

620 C.3 Supporting lemmas

621 **Lemma C.2.** Suppose Assumption 2.2 holds. With probability $1 - \gamma$ for some failure probability
622 $\gamma > 0$, we have

$$\sum_{k=1}^K \sum_{h=1}^H \frac{1}{\sqrt{N_h^k(s_h^k, a_h^k)}} \leq \left\lceil \frac{\log \frac{HK}{\gamma}}{-\log(1 - \lambda_0^2)} \right\rceil \sqrt{S A K H}.$$

623 *Proof of Lemma C.2.* For any time h , we denote $\mathcal{K}^e(h)$ as the collection of episodes that the h -th

624 and $(h+1)$ -th step observations are available. It is clear that the cardinality of $\mathcal{K}^e(h)$ is bounded

625 by K for any h . Within each $\mathcal{K}^e(h)$, we would like to bound the gap between two observations.

626 Thanks to Lemma C.3, the gap is bounded by q with probability $1 - K(1 - \lambda_0^2)^{q+1}$. We set

627 $K(1 - \lambda_0^2)^{q+1} = \gamma/H$, which implies $q = \left\lceil \frac{\log \frac{HK}{\gamma}}{-\log(1 - \lambda_0^2)} \right\rceil$. Therefore, for any time step h , available

628 observations are at most separated by q episodes.

629 With these notations, we bound

$$\begin{aligned} \sum_{k=1}^K \sum_{h=1}^H \frac{1}{\sqrt{N_h^k(s_h^k, a_h^k)}} &\stackrel{(i)}{\leq} \left\lceil \frac{\log \frac{HK}{\gamma}}{-\log(1 - \lambda_0^2)} \right\rceil \sum_{h=1}^H \sum_{k \in \mathcal{K}^{\text{eff}}(h)} \frac{1}{\sqrt{N_h^k(s_h^k, a_h^k)}} \\ &\stackrel{(ii)}{\leq} \left\lceil \frac{\log \frac{HK}{\gamma}}{-\log(1 - \lambda_0^2)} \right\rceil \sum_{h=1}^H \sum_{k=1}^K \frac{1}{\sqrt{N_h^k(s_h^k, a_h^k)}} \\ &\stackrel{(iii)}{\leq} 2 \left\lceil \frac{\log \frac{HK}{\gamma}}{-\log(1 - \lambda_0^2)} \right\rceil \sqrt{S A H K}, \end{aligned}$$

630 where inequality (i) follows since N_h^k will only be updated when $h \in \mathcal{K}^e$ (h) and then repeat at
631 most $\left\lceil \frac{\log \frac{HK}{\log(1-\lambda_0^2)}}{\log(1-\lambda_0^2)} \right\rceil$ times, inequality (ii) invokes the cardinality bound of \mathcal{K}^e (h), and inequality
632 (iii) follows from the standard pigeon-hole principle. \square

633 **Lemma C.3.** Let $\{u_i\}_{i=1}^k$ be i.i.d. Bernoulli random variables. Suppose $\mathbb{P}(u_i = 1) = \theta$. Define the
634 largest gap between u_i 's as

$$g(k) = \sup\{j - i : u_i = 0 \text{ and } u_j = 0 \text{ with } u_\ell = 1 \text{ for } \ell = i + 1, \dots, j - 1\}.$$

635 Then for any integer $q > 0$, the following tail probability bound holds

$$\mathbb{P}(g(k) > q) \leq k\theta^{q+1}.$$

636 *Proof of Lemma C.3.* We denote $I_{\text{neg}} = \{\ell_1, \dots, \ell_m\}$ as the index set for $u_{\ell_i} = 0$ when $i =$
637 $1, \dots, |I_{\text{neg}}|$. Let $v_j = \ell_{j+1} - \ell_j$, which is a geometric random variable with a success rate θ . Note
638 that the cardinality of I_{neg} is at most k . Therefore, we have

$$\begin{aligned} \mathbb{P}(g(k) > q) &\leq \mathbb{P}(\max_{j=1, \dots, k} v_j > q) \\ &= 1 - \mathbb{P}(v_j \leq q \text{ for } j = 1, \dots, k) \\ &= 1 - (1 - \theta^{q+1})^k \\ &\leq k\theta^{q+1}, \end{aligned}$$

639 where the last inequality follows from $1 - k\theta^{q+1} \leq (1 - \theta^{q+1})^k$. \square

640 D Helper concentration inequalities

641 **Lemma D.1** (Bernstein's inequality). Let x_1, \dots, x_n be i.i.d. zero mean random variables. Suppose
642 $|x_i| \leq M$ for any $i = 1, \dots, n$. Then for all positive t , it holds

$$\mathbb{P}\left(\sum_{i=1}^n x_i > t\right) \leq \exp\left(-\frac{\frac{1}{2}t^2}{\sum_{i=1}^n \text{Var}[x_i] + \frac{1}{3}Mt}\right).$$

643 In particular, given a failure probability $\gamma < 1$, it holds

$$\mathbb{P}\left(\sum_{i=1}^n x_i > \sqrt{\sum_{i=1}^n \text{Var}[x_i] \log \frac{1}{\gamma}} + M \log \frac{1}{\gamma}\right) \leq \gamma.$$

644 *Proof of Lemma D.1.* The proof of Bernstein's inequality is standard; see for example [Wainwright](#)
645 [\[2019, Section 2.1\]](#). Here we verify the second claim. Let $\exp\left(-\frac{\frac{1}{2}t^2}{\sum_{i=1}^n \text{Var}[x_i] + \frac{1}{3}Mt}\right) \leq \gamma$ hold true.
646 We find a suitable t by

$$\begin{aligned} &\exp\left(-\frac{\frac{1}{2}t^2}{\sum_{i=1}^n \text{Var}[x_i] + \frac{1}{3}Mt}\right) \leq \gamma \\ \iff &\frac{\frac{1}{2}t^2}{\sum_{i=1}^n \text{Var}[x_i] + \frac{1}{3}Mt} \geq \log \frac{1}{\gamma} \\ \iff &t^2 - \frac{2}{3}tM \log \frac{1}{\gamma} \geq \sum_{i=1}^n \text{Var}[x_i] \log \frac{1}{\gamma} \\ \iff &t \geq \sqrt{\sum_{i=1}^n \text{Var}[x_i] \log \frac{1}{\gamma} + \frac{1}{9}M^2 \log^2 \frac{1}{\gamma} + \frac{1}{3}M \log \frac{1}{\gamma}}. \end{aligned}$$

647 It is enough to choose $t = \sqrt{\sum_{i=1}^n \text{Var}[x_i] \log \frac{1}{\gamma} + M \log \frac{1}{\gamma}}$. \square

648 **Lemma D.2** (Hoeffding's inequality). Let x_1, \dots, x_n be i.i.d. random variables. Suppose $a_i \leq x_i \leq$
649 b_i for any $i = 1, \dots, n$. Then for all positive t , it holds

$$\mathbb{P} \left(\left| \sum_{i=1}^n x_i - \mathbb{E} \left[\sum_{i=1}^n x_i \right] \right| > t \right) \leq 2 \exp \left(- \frac{2t^2}{\sum_{i=1}^n (b_i - a_i)^2} \right).$$

650 In particular, given a failure probability $\gamma < 1$, it holds

$$\mathbb{P} \left(\frac{1}{n} \left| \sum_{i=1}^n x_i - \mathbb{E} \left[\sum_{i=1}^n x_i \right] \right| > \sqrt{\frac{\sum_{i=1}^n (b_i - a_i)^2 \log \frac{2}{\gamma}}{2n^2}} \right) \leq \gamma.$$

651 *Proof of Lemma D.2.* The proof is standard; see [Wainwright \[2019, Section 2.1\]](#). □

652 **Lemma D.3** (Azuma-Hoeffding's inequality). Let x_1, \dots, x_n be a martingale adapted to filtration
653 $\mathcal{F}_1 \subset \dots \subset \mathcal{F}_n$. Suppose $\mathbb{E}[x_i - \mathbb{E}[x_i | \mathcal{F}_{i-1}]] = 0$ and $|x_i - \mathbb{E}[x_i]| \leq c_i$. Then for all positive t , it
654 holds

$$\mathbb{P} \left(\sum_{i=1}^n x_i - \mathbb{E}[x_i] > t \right) \leq \exp \left(- \frac{t^2}{2 \sum_{i=1}^n c_i^2} \right).$$

655 In particular, given a failure probability $\gamma < 1$, it holds

$$\mathbb{P} \left(\sum_{i=1}^n x_i - \mathbb{E}[x_i] > \sqrt{2 \sum_{i=1}^n c_i^2 \log \frac{1}{\gamma}} \right) \leq \gamma.$$

656 *Proof of Lemma D.3.* The proof is standard and applies [Lemma D.2](#). □