
Preference Elicitation for Offline Reinforcement Learning

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Applying reinforcement learning (RL) to real-world problems is often made chal-
2 lenging by the inability to interact with the environment and the difficulty of de-
3 signing reward functions. Offline RL addresses the first challenge by consider-
4 ing access to an offline dataset of environment interactions labeled by the reward
5 function. In contrast, Preference-based RL does not assume access to the reward
6 function and learns it from preferences, but typically requires an online interaction
7 with the environment. We bridge the gap between these frameworks by explor-
8 ing efficient methods for acquiring preference feedback in a fully offline setup.
9 We propose Sim-OPRL, an offline preference-based reinforcement learning algo-
10 rithm, which leverages a learned environment model to elicit preference feed-
11 back on simulated rollouts. Drawing on insights from both the offline RL and
12 the preference-based RL literature, our algorithm employs a pessimistic approach
13 for out-of-distribution data, and an optimistic approach for acquiring informative
14 preferences about the optimal policy. We provide theoretical guarantees regarding
15 the sample complexity of our approach, dependent on how well the offline data
16 covers the optimal policy. Finally, we demonstrate the empirical performance of
17 Sim-OPRL in different environments.

18 1 Introduction

19 While reinforcement learning (RL) [Sutton and Barto, 2018] achieves excellent performance in var-
20 ious decision-making tasks [Mirhoseini et al., 2020, Degraeve et al., 2022], its practical deployment
21 remains limited by the requirement of direct interaction with the environment. This can be imprac-
22 tical or unsafe in real-world scenarios. For example, patient management and treatment in intensive
23 care units involve complex decision-making that has often been framed as a reinforcement learning
24 problem [Raghu et al., 2017]. However, the timing, dosage, and combination of treatments required
25 are critical to patient safety, and incorrect decisions can lead to severe complications or death, mak-
26 ing the use of traditional RL algorithms unfeasible [Tang and Wiens, 2021]. Offline RL emerges as
27 a promising solution, allowing policy learning from entirely observational data [Levine et al., 2020].

28 Still, a challenge with Offline RL is its requirement for an explicit reward function. Quantifying
29 the numerical value of taking a certain action in a given environment state is often challenging [Yu
30 et al., 2021]. Preference-based RL offers a compelling alternative, relying on comparisons between
31 different trajectories, and being often easier for humans to provide [Wirth et al., 2017]. In medical
32 settings, for instance, clinicians may be queried for feedback on which trajectories lead to favorable
33 outcomes. Unfortunately, most algorithms for preference acquisition require environment interac-
34 tion [Saha et al., 2023, Chen et al., 2022] and are therefore not applicable to the offline setting.

35 Lack of environment interaction and reward learning are thus two critical challenges for real-world
36 RL deployment that are rarely tackled jointly. In this work, we address the problem of prefer-
37 ence elicitation for offline reinforcement learning by asking: *What trajectories should we sample*

38 *to minimize the number of human queries required to learn the best offline policy?* This presents a
 39 challenging problem as it combines learning from offline data and active feedback acquisition, two
 40 frameworks that require opposing inductive biases for conservatism and exploration, respectively.

41 To the best of our knowledge, the only strategy proposed in prior work is to acquire feedback directly
 42 over samples within an offline dataset of trajectories [Shin et al., 2022, Offline Preference-based Re-
 43 ward Learning (OPRL)]. We propose an alternative solution that queries feedback on *simulated*
 44 *rollouts* by leveraging a learned environment model. Our offline preference-based reinforcement
 45 learning algorithm, Sim-OPRL, strikes a balance between conservatism and exploration by combin-
 46 ing pessimism when handling states out-of-distribution from the observational data [Jin et al., 2021,
 47 Zhan et al., 2023a], and optimism in acquiring informative preferences about the optimal policy
 48 [Saha et al., 2023, Chen et al., 2022]. We validate our approach through both theoretical and empir-
 49 ical analysis, demonstrating the superior performance of Sim-OPRL across various environments.

50 Our contributions are the following: (1) In Section 3, we first formalize the new problem setting of
 51 preference elicitation for offline reinforcement learning, which allows for **complementing offline**
 52 **data with preference feedback**. This framework is crucial for real-world applications where direct
 53 environment interaction is infeasible and reward functions are challenging to design manually, yet
 54 experts can be queried for their knowledge. (2) In Section 4, we propose a novel offline preference-
 55 based RL algorithm that is independent of the specific preference elicitation strategy and recovers
 56 a robust policy from an offline dataset and preference feedback. (3) Next, in Section 5, we provide
 57 theoretical guarantees on eliciting preferences over samples from the offline dataset, complementing
 58 work from Shin et al. [2022]. (4) Then, in Section 6 we propose our own **efficient preference**
 59 **elicitation algorithm** based on simulated rollouts in a learned environment model. (5) Finally, we
 60 establish the **theoretical guarantees** of our algorithm and demonstrate its **empirical efficiency** and
 61 scalability in different decision-making environments.

62 2 Related Work

63 Our problem setting shares similarities with Offline RL and Preference-based RL, which we sum-
 64 marize below. We position ourselves relative to our closest related works in Table 1.

65 **Offline RL.** Offline Reinforcement Learning has gained significant traction in recent years, as the
 66 practicality of training RL agents without environment interaction makes it relevant to real-world
 67 applications [Levine et al., 2020]. However, learning from observational data only is a source of
 68 bias in the model, as the data may not cover the entire state-action space. Offline RL algorithms
 69 therefore output pessimistic policies, which has been shown to minimize suboptimality Jin et al.
 70 [2021]. Model-based approaches show particular promise for their sample efficiency [Yu et al.,
 71 2020, Kidambi et al., 2020, Uehara and Sun, 2021]. In this work, we study the setting where reward
 72 signals are unavailable and must be estimated by actively querying preference feedback.

73 **Preference-based RL.** Rather than accessing numerical reward values for each state-action pair as
 74 in traditional online RL, preference-based RL learns the reward model through collecting pairwise
 75 preferences over trajectories [Wirth et al., 2017]. Different preference elicitation strategies have
 76 been proposed for this framework, generally based on knowing the transition model exactly or on
 77 having access to the environment for rollouts [Christiano et al., 2017, Saha et al., 2023, Chen et al.,
 78 2022, Lindner et al., 2021, Zhan et al., 2023b, Sadigh et al., 2018, Brown et al., 2020].

79 **Offline Preference-based RL.** The development of preference-based RL algorithms based on of-
 80 fline data only is critical to settings where environment interaction is not feasible for safety and
 81 efficiency reasons. Still, this framework remains largely unexplored in the literature. While Zhu
 82 et al. [2023], Zhan et al. [2023a] demonstrate the value of pessimism in offline preference-based re-
 83 inforcement learning, they do not consider how to query feedback actively. On the other hand, Shin

Table 1: Comparison of related work on preference elicitation.

Framework	Offline	Efficient Sampling	Robustness Guarantees	Practical Implementation
PbOP [Chen et al., 2022]	✗	✓	✓	✗
REGIME [Zhan et al., 2023b]	✗	✓	✓	✗
FREEHAND [Zhan et al., 2023a]	✓	✗	✓	✗
OPRL [Shin et al., 2022]	✓	✓ ¹	✗	✓
Sim-OPRL (Ours)	✓	✓	✓	✓

¹We demonstrate this in the present work.

84 et al. [2022] propose an empirical comparison of different preference sampling trajectories from an
 85 offline trajectories buffer. In Section 5, we provide a theoretical analysis of their approach, then
 86 propose an alternative sampling strategy based on simulated trajectory rollouts in Section 6, which
 87 enjoys both theoretical and empirical motivation.

88 3 Problem formulation

89 3.1 Preliminaries

90 **Markov Decision Process.** We consider the episodic Markov Decision Process (MDP), defined by
 91 the tuple $\mathcal{M} = (\mathcal{S}, \mathcal{A}, H, T, R)$, where \mathcal{S} is the state space, \mathcal{A} is the action space, H is the episode
 92 length, $T : \mathcal{S} \times \mathcal{A} \rightarrow \Delta_{\mathcal{S}}$ is the transition function, $R : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ is the reward function.
 93 We assume an initial state s_0 , but our analysis could be easily generalized to a fixed initial state
 94 distribution. At time t , the environment is at state $s_t \in \mathcal{S}$ and an agent selects an action $a_t \in \mathcal{A}$. The
 95 agent then receives a reward $R(s_t, a_t)$ and the environment transitions to state $s_{t+1} \sim T(\cdot | s_t, a_t)$.
 96 We describe an agent’s behavior through a policy function $\pi : \mathcal{S} \rightarrow \Delta_{\mathcal{A}}$, such that $\pi(a|s)$ is the
 97 probability of taking action a in state s . Let $\tau = (s_0, a_0, \dots, s_H, a_H)$ denote the trajectory of state-
 98 action pairs of an interaction episode with the environment. With an abuse of notation, we also write
 99 $R(\tau) = \sum_t R(s_t, a_t)$. Let d_T^π denote the distribution of trajectories induced by rolling out policy
 100 π in transition model T . We denote the expected return of policy π as $V_{T,R}^\pi = \mathbb{E}_{\tau \sim d_T^\pi}[R(\tau)]$, and
 101 $\pi^* = \operatorname{argmax}_\pi V_{T,R}^\pi$ denotes the optimal policy in \mathcal{M} .

102 **Preference-based Reinforcement Learning.** Rather than observing numerical rewards at each state
 103 and action, we receive preference feedback over trajectories. For a pair of trajectories (τ_1, τ_2) , we
 104 obtain binary feedback $o \in \{0, 1\}$ about whether τ_1 is preferred to τ_2 . We assume that preference
 105 labels follow the Bradley-Terry model [Bradley and Terry, 1952]:

$$P_R(\tau_1 \succ \tau_2) := P(o = 1 | \tau_1, \tau_2) = \frac{\exp(R(\tau_1))}{\exp(R(\tau_1)) + \exp(R(\tau_2))} = \sigma(R(\tau_1) - R(\tau_2)), \quad (1)$$

106 where \succ denotes a preference relationship and σ is the sigmoid function. Within this framework,
 107 *preference elicitation* refers to the process of sampling preferences to obtain information about both
 108 the preference function and the system dynamics [Wirth et al., 2017].

109 3.2 Offline Preference Elicitation

110 We assume access to an observational dataset of trajectories $\mathcal{D}_{offline} = \{\tau : \tau \sim d_T^{\pi_\beta}\}$, where π_β
 111 is an unknown behavioural policy in \mathcal{M} . As in Offline RL, we do *not* have access to the decision-
 112 making environment to observe transition dynamics or rewards under alternative action choices. We
 113 assume *not* to have access to the reward function, but we can query preference feedback from a
 114 human to generate a dataset of preferences $\mathcal{D}_{pref} = \{(\tau_1, \tau_2, o)\}$.

115 **Optimality Criterion.** Based only on our offline dataset $\mathcal{D}_{offline}$, our goal is to recover a policy
 116 $\hat{\pi}^*$ that minimizes suboptimality in the true environment with as few human preference queries
 117 as possible. Let $\pi_{offline}^*$ denote the *optimal offline policy* estimated based on the offline data, with
 118 access to the true reward function R , and let ϵ_T denote its suboptimality. Since preference elicitation
 119 only allows us to estimate the reward function, we do not aim to achieve a suboptimality less than
 120 ϵ_T .² Our objective is then formalized as follows.

121 **Definition 3.1** (Optimality Criterion of Offline Preference Elicitation). *Let π^* be the optimal policy*
 122 *in \mathcal{M} and $\hat{\pi}^*$ be the estimated optimal policy based on an offline dataset $\mathcal{D}_{offline}$ and $N_p > 0$*
 123 *preference queries. Let ϵ_T be the inherent suboptimality assuming access to the true reward function.*
 124 *We say that a sampling strategy is (ϵ, δ, N_p) -correct if for every $\epsilon \geq \epsilon_T$, with probability at least*
 125 *$(1 - \delta)$, it holds that $V_{T,R}^{\pi^*} - V_{T,R}^{\hat{\pi}^*} \leq \epsilon$.*

126 Our work is the first to formalize this important problem, which faces the challenge of balancing
 127 exploration when actively acquiring feedback and bias mitigation in learning from offline data.

128 **Function classes.** We estimate the reward function and transition kernel with general function ap-
 129 proximation; let \mathcal{F}_R and \mathcal{F}_T denote the classes of functions considered respectively. We also assume
 130 a policy class Π . Our theoretical analysis also requires the following assumptions and definitions,
 131 which are standard in preference-based RL [Chen et al., 2022, Zhan et al., 2023a].

²However, ϵ_T is not formally a lower bound for our problem, as shown in Appendix A.3.

Algorithm 1 Offline Preference-based Reinforcement Learning with Preference Elicitation

Input: Observational trajectories dataset $\mathcal{D}_{offline}$. Significance $\delta \in (0, 1)$, preference budget N_p .

Output: $\hat{\pi}^*$

- 1: Estimate \hat{T} and u_T via maximum likelihood over the observational data $\mathcal{D}_{offline}$.
 - 2: $\mathcal{D}_{pref} \leftarrow \emptyset$.
 - 3: **for** $k = 1, \dots, N_p$ **do**
 - 4: Generate trajectory pairs (τ_1, τ_2) . ▷ **Preference Elicitation:** Sections 5 and 6
 - 5: Collect preference label o for (τ_1, τ_2) .
 - 6: $\mathcal{D}_{pref} \leftarrow \mathcal{D}_{pref} \cup \{(\tau_1, \tau_2, o)\}$.
 - 7: Estimate \hat{R} and u_R via maximum likelihood over the preference data \mathcal{D}_{pref} .
 - 8: **end for**
 - 9: $\hat{\pi}^* \leftarrow \operatorname{argmax}_{\pi \in \Pi} \mathbb{E}_{\tau \sim d_T^\pi} [\hat{R}(\tau) - u_R(\tau) - u_T(\tau)]$
-

132 **Assumption 3.1** (Realizability). *The true reward function belongs to the reward class: $R \in \mathcal{F}_R$.*
 133 *The true transition function belongs to the transition class: $T \in \mathcal{F}_T$. The optimal policy belongs to*
 134 *the policy class: $\pi^* \in \Pi$.*

135 **Assumption 3.2** (Boundedness). *The reward function is bounded: $0 \leq \tilde{R}(\tau) \leq R_{max}$ for all*
 136 *$\tilde{R} \in \mathcal{F}_R$ and all trajectories τ .*

137 **Definition 3.2** (ϵ -bracketing number). *Let \mathcal{F} be a class of real functions $f : \mathcal{X} \rightarrow \mathbb{R}$. We say (l, u)*
 138 *is an ϵ -bracket if $l(x) \leq u(x)$ and $\|u(x) - l(x)\|_1 \leq \epsilon$ for all $x \in \mathcal{X}$. The ϵ -bracketing number of*
 139 *\mathcal{F} , denoted $\mathcal{N}_{\mathcal{F}}(\epsilon)$, is the minimal number of ϵ -brackets $(l^n, u^n)_{n=1}^N$ needed so that for any $f \in \mathcal{F}$,*
 140 *there is a bracket $i \in [N]$ containing it, meaning $l^i(x) \leq f(x) \leq u^i(x)$ for all $x \in \mathcal{X}$.*

141 Let $\mathcal{N}_{\mathcal{F}_R}(\epsilon)$ and $\mathcal{N}_{\mathcal{F}_T}(\epsilon)$ denote the ϵ -bracketing numbers of \mathcal{F}_R and \mathcal{F}_T respectively. This measures
 142 the complexity of the function classes [Geer, 2000].

143 **Definition 3.3** (Transition concentrability coefficient, Zhan et al. [2023a]). *The concentrability co-*
 144 *efficient w.r.t. transition classes \mathcal{F}_T and the optimal policy π^* is defined as:*

$$C_T(\mathcal{F}_T, \pi^*) = \sup_{\tilde{T} \in \mathcal{F}_T} \left[\frac{\mathbb{E}_{(s,a) \sim d_T^{\pi^*}} [|T(\cdot|s, a) - \tilde{T}(\cdot|s, a)|]}{\sqrt{\mathbb{E}_{(s,a) \sim \mathcal{D}_{offline}} [|T(\cdot|s, a) - \tilde{T}(\cdot|s, a)|^2]}} \right]$$

145

146 The concentrability coefficient measures the coverage of the optimal policy in the offline
 147 dataset. Note that C_T is upper-bounded by the density-ratio coefficient: $C_T(\mathcal{F}_T, \pi^*) \leq$
 148 $\sup_{(s,a) \in \mathcal{S} \times \mathcal{A}} d_T^{\pi^*}(s, a) / d_T^{\pi_\beta}(s, a)$, where π_β is the behavioural policy underlying $\mathcal{D}_{offline}$.

149 4 Offline Preference-based RL with Preference Elicitation

150 In this section, we propose a general framework for offline preference-based reinforcement learning.
 151 The next two sections propose two different preference elicitation strategies. As learning must be
 152 carried out in two stages, with environment dynamics based on $\mathcal{D}_{offline}$ and reward learning on
 153 \mathcal{D}_{pref} , we adopt a model-based approach which we summarize in Algorithm 1.

154 **Model Learning.** We first leverage the offline data to learn a model of the environment dynamics,
 155 fitting a transition model \hat{T} and an uncertainty function u_T through maximum likelihood:

$$\begin{aligned} \hat{T} &= \operatorname{argmax}_{\tilde{T} \in \mathcal{F}_T} \mathbb{E}_{(s,a,s') \sim \mathcal{D}_{offline}} \left[\log \tilde{T}(s'|s, a) \right], \\ u_T(s, a) &= \max_{\tilde{T}_1, \tilde{T}_2 \in \mathcal{T}} |\tilde{T}_1(\cdot|s, a) - \tilde{T}_2(\cdot|s, a)| \cdot R_{max}, \end{aligned}$$

156 where $\mathcal{T} = \{\tilde{T} \in \mathcal{F}_T \mid \mathbb{E}_{(s,a,s') \sim \mathcal{D}_{offline}} \left[\log \hat{T}(s'|s, a) / \tilde{T}(s'|s, a) \right] \leq \beta_T\}$, defining a confidence
 157 set over the MLE estimate, and β_T is a hyperparameter. In a practical implementation, this can be
 158 achieved by training an ensemble of models on different data bootstraps [Lakshminarayanan et al.,
 159 2017].

160 **Iterative Preference Elicitation and Reward Learning.** As with the transition model, our algo-
 161 rithm estimates the reward function \hat{R} and its uncertainty function through maximum likelihood

162 over iteratively collected preference data \mathcal{D}_{pref} :

$$\begin{aligned} \hat{R} &= \operatorname{argmax}_{\tilde{R} \in \mathcal{F}_R} \mathbb{E}_{(\tau_1, \tau_2, o) \sim \mathcal{D}_{pref}} [o \log P_{\tilde{R}}(\tau_1 \succ \tau_2) + (1 - o) \log P_{\tilde{R}}(\tau_2 \succ \tau_1)], \\ u_R(\tau) &= \max_{\tilde{R}_1, \tilde{R}_2 \in \mathcal{R}} |\tilde{R}_1(\tau) - \tilde{R}_2(\tau)|, \end{aligned}$$

163 where $\mathcal{R} = \{\tilde{R} \in \mathcal{F}_R \mid \mathbb{E}_{(\tau_1, \tau_2, o) \sim \mathcal{D}_{pref}} [\log P_{\tilde{R}}(\tau_1 \succ \tau_2) / P_{\tilde{R}}(\tau_1 \succ \tau_2)] \leq \beta_R\}$ defines the confi-
164 dence set and β_R is a hyperparameter. We also define preference uncertainty between two trajec-
165 tories τ_1, τ_2 as $u_{P_R}(\tau_1, \tau_2) = \max_{\tilde{R}_1, \tilde{R}_2 \in \mathcal{R}} |P_{\tilde{R}_1}(\tau_1 \succ \tau_2) - P_{\tilde{R}_2}(\tau_1 \succ \tau_2)|$.

166 The choice of trajectory sampling strategy for preference elicitation in line 4 is critical to efficiently
167 obtaining an ϵ -optimal policy. We present two possible strategies in Sections 5 and 6. Note that by
168 focusing on sample efficiency as in prior work on preference elicitation [Chen et al., 2022], we do not
169 necessarily optimize for computational efficiency; this could be improved by collecting preferences
170 in batches to reduce the number of reward training loops.

171 **Pessimistic Policy Optimization.** Finally, our algorithm outputs a policy $\hat{\pi}^*$ that is optimal while
172 ensuring robustness to modeling error. This means optimizing for the worst-case value function over
173 the remaining transition and reward uncertainties [Levine et al., 2020]:

$$\hat{\pi}^* = \operatorname{argmax}_{\pi \in \Pi} \min_{\tilde{T} \in \mathcal{T}, \tilde{R} \in \mathcal{R}} V_{\tilde{T}, \tilde{R}}^{\pi}.$$

174 This analysis provides a worst-case robustness guarantee when considering well-calibrated confi-
175 dence intervals, as detailed in Sections 5.1 and 6.1. For a practical implementation of our algorithm,
176 we penalize the reward function by the uncertainty as in model-based offline RL methods [Yu et al.,
177 2020, Chang et al., 2021]. Our optimal robust policy therefore maximizes the following objective:

$$\hat{\pi}^* = \operatorname{argmax}_{\pi \in \Pi} \mathbb{E}_{\tau \sim d_{\tilde{T}}^{\pi}} [\hat{R}(\tau) - u_R(\tau) - u_T(\tau)]. \quad (2)$$

178 We show in Appendix A.2 that this is indeed a lower bound of the true value function. This objective
179 allows for controlling the degree of conservatism in practice through the width of the confidence
180 intervals used to determine u_R and u_T .

181 5 Preference Elicitation from Offline Trajectories

182 A first strategy for preference elicitation without environment interaction is to sample trajectories
183 directly from the offline dataset. Shin et al. [2022] propose this approach as Offline Preference-based
184 Reward Learning (OPRL), and design a uniform and uncertainty-sampling variant:

$$\begin{aligned} \text{OPRL Uniform Sampling:} & \quad \tau_1, \tau_2 \sim \mathcal{D}_{offline} \\ \text{OPRL Uncertainty Sampling:} & \quad \tau_1, \tau_2 = \operatorname{argmax}_{\tau_1, \tau_2 \in \mathcal{D}_{offline}} u_{P_R}(\tau_1, \tau_2) \end{aligned}$$

185 5.1 Theoretical Guarantees.

187 We obtain the following result, demonstrated in Appendix A.4. The suboptimality of the estimated
188 policy $\hat{\pi}^*$ is bounded by the policy evaluation error for the optimal policy π^* . This error decomposes
189 into a term depending on transition model estimation, and one on reward model estimation.

190 **Theorem 5.1.** *For any $\delta \in (0, 1]$, let $\beta_T = c'_T \log(HN_{\mathcal{F}_T}(1/N_o)/\delta)/N_o$ and $\beta_R =$
191 $c'_R \log(\mathcal{N}_{\mathcal{F}_R}(1/N_p)/\delta)/N_p$, where $N_o = H|\mathcal{D}_{offline}|$ is the number of observed transitions in the
192 observational dataset and c'_T, c'_R are universal constants. The policy $\hat{\pi}^*$ estimated by Algorithm 1,
193 with preference elicitation based on offline trajectories, achieves the following suboptimality with
194 probability $1 - \delta$:*

$$V^{\pi^*} - V^{\hat{\pi}^*} \leq \underbrace{HR_{max} C_T(\mathcal{F}_T, \pi^*) \sqrt{\frac{c_T \log(HN_{\mathcal{F}_T}(1/N_o)/\delta)}{N_o}}}_{\text{transition term } \epsilon_T} + 2\alpha\kappa C_R(\mathcal{F}_R, \pi^*) \underbrace{\sqrt{\frac{c_R \log(\mathcal{N}_{\mathcal{F}_R}(1/N_p)/\delta)}{N_p}}}_{\text{reward term}},$$

195 where $\alpha = 1$ for uniform sampling or $\alpha \leq 1$ for uncertainty sampling, C_R is a concentrability
196 measure for the reward function, $\kappa = \sup_{r \in [-R_{max}, R_{max}]} \frac{1}{\sigma'(r)}$ measures the degree of non-linearity
197 of the sigmoid function, and c_T, c_R are universal constants.

198 In the special case where both the transition and reward functions are learned on a fixed initial
199 preference dataset (no preference elicitation), we recover Theorem 1 from Zhan et al. [2023a]. Im-
200 portantly, parameter α allows us to motivate the superior efficiency of uncertainty sampling over
201 uniform sampling, observed empirically in Shin et al. [2022] and in Section 7.

202 6 Preference Elicitation from Simulated Trajectories

203 We now propose our alternative strategy for generating trajectories for offline preference elicitation:
 204 **Simulated Offline Preference-based Reward Learning (Sim-OPRL)**. This method simulates tra-
 205 jectories (τ_1, τ_2) by leveraging the *learned environment model*. This overcomes a limitation of
 206 OPRL, which is only designed to reduce uncertainty about the reward functions in \mathcal{R} , by instead
 207 reducing uncertainty about which policies are plausibly optimal. Our approach is inspired by effi-
 208 cient online preference elicitation algorithms [Saha et al., 2023, Chen et al., 2022], which we modify
 209 for practical implementation. We account for the offline nature of our problem by avoiding regions
 210 out of the distribution of the data: the sampling strategy is optimistic with respect to uncertainty in
 211 rewards, but pessimistic with respect to uncertainty in transitions.

212 We summarize our approach to generating simulated trajectories for preference elicitation in Al-
 213 gorithm 2 and refer the reader to Appendix B for practical implementation details. First, we con-
 214 struct a set of **candidate optimal policies** $\Pi_{offline}$, containing policy $\pi_{offline}^*$ (optimal under the
 215 pessimistic model and the true reward function) with high probability – as demonstrated in Ap-
 216 pendix A.5.2. Next, within this set of candidate policies, we identify the two most **exploratory**
 217 **policies** π_1, π_2 , chosen to maximize preference uncertainty u_{PR} . Finally, we roll out these policies
 218 within our learned transition model to generate a trajectory pair (τ_1, τ_2) for preference feedback.

Algorithm 2 Preference Elicitation through Simulated Trajectory Sampling.

Input: Pessimistic transition model \hat{T}_{inf} . Reward confidence set \mathcal{R} and preference uncertainty function u_{PR} .

Output: (τ_1, τ_2)

- 1: Estimate optimal offline policy set: $\Pi_{offline} = \{\pi \mid \pi = \operatorname{argmax}_{\pi \in \Pi} \mathbb{E}_{\tau \sim d_{\hat{T}_{inf}}^{\pi}} [\tilde{R}(\tau)] \forall \tilde{R} \in \mathcal{R}\}$.
 - 2: Identify exploratory policies: $\pi_1, \pi_2 = \operatorname{argmax}_{\pi_1, \pi_2 \in \Pi_{offline}} \mathbb{E}_{\tau_1 \sim d_{\hat{T}_{inf}}^{\pi_1}, \tau_2 \sim d_{\hat{T}_{inf}}^{\pi_2}} [u_{PR}(\tau_1, \tau_2)]$
 - 3: Rollouts in model: $\tau_1 \sim d_{\hat{T}_{inf}}^{\pi_1}, \tau_2 \sim d_{\hat{T}_{inf}}^{\pi_2}$.
-

219 6.1 Theoretical Guarantees

220 We decompose suboptimality in a similar way to Section 5.1, but obtain a reward suboptimality term
 221 that depends on the learned dynamics model instead of the true one, and on $\pi_{offline}^*$ instead of π^* :

$$V^{\pi^*} - V^{\hat{\pi}^*} \leq \underbrace{(V_{T,R}^{\pi^*} - V_{\hat{T}_{inf},R}^{\pi^*})}_{\text{transition term } \epsilon_T} + \underbrace{(V_{\hat{T}_{inf},R}^{\pi_{offline}^*} - V_{\hat{T}_{inf},\hat{R}_{inf}}^{\pi_{offline}^*})}_{\text{reward term}}. \quad (3)$$

222 Analysis of the suboptimality due to transition error is identical to above, but the reward term is
 223 thus significantly different. By design, our sampling strategy ensures good coverage of preferences
 224 over $\pi_{offline}^*$ within the learned environment model, which **eliminates the concentrability term**
 225 **for the reward** C_R . We refer the reader to Appendix A.5 for the proof of Theorem 6.1.

226 **Theorem 6.1.** *For any $\delta \in (0, 1]$, let $\beta_T = c'_T \log(HN_{\mathcal{F}_T}(1/N_o)/\delta)/N_o$ and $\beta_R =$
 227 $c'_R \log(\mathcal{N}_{\mathcal{F}_R}(1/N_p)/\delta)/N_p$, where $N_o = H|\mathcal{D}_{offline}|$ is the number of observed transitions in the
 228 observational dataset and c'_T, c'_R are universal constants. The policy $\hat{\pi}^*$ estimated by Algorithm 1,
 229 with a preference sampling strategy based on rollouts in the learned transition model, achieves the
 230 following suboptimality with probability $1 - \delta$:*

$$V^{\pi^*} - V^{\hat{\pi}^*} \leq HR_{max} C_T(\mathcal{F}_T, \pi^*) \sqrt{\frac{c_T \log(HN_{\mathcal{F}_T}(1/N_o)/\delta)}{N_o}} + 2\kappa \sqrt{\frac{c_R \log(\mathcal{N}_{\mathcal{F}_R}(1/N_p)/\delta)}{N_p}}.$$

231 6.2 Discussion

232 Our theoretical results demonstrate that the learned policy can achieve performance comparable to
 233 the optimal policy, and thus satisfy our optimality criterion in Definition 3.1, provided it is covered
 234 by the offline data ($C_T(\mathcal{F}_T, \pi^*), C_R(\mathcal{F}_R, \pi^*) < \infty$). Empirical results in Section 7 confirm that
 235 performance is poor when the behavioral policy is suboptimal, inducing a large C_T or C_R .

236 **Offline Trajectories vs. Simulated Rollouts.** While both OPRL and Sim-OPRL depend on the
 237 offline dataset for estimating environment dynamics, they induce different suboptimality in model-
 238 ing preference feedback. Simulated rollouts are designed to achieve good coverage of the optimal

239 offline policy $\pi_{offline}^*$, which avoids wasting preference budget on trajectories with low rewards or
 240 high transition uncertainty. In contrast, as shown in Zhan et al. [2023a], due to the dependence of
 241 preferences on full trajectories, the reward concentrability term C_R in Theorem 5.1 can be large.

242 **Transition vs. Preference Model Quality.** Our theoretical analysis also suggests an interesting
 243 trade-off in the sample efficiency of our approach, depending on the accuracy of the transition model.
 244 The width of the confidence interval reduces as significance parameter δ or dataset size increase,
 245 or as function class complexity $\mathcal{N}_{\mathcal{F}_T}$ decreases. For a target suboptimality gap ϵ , provided the
 246 optimal offline policy $\pi_{offline}^*$ has a gap $\epsilon_T < \epsilon$, then the number of preferences required is of the
 247 order of $\mathcal{O}(\log(1/\delta)/(\epsilon - \epsilon_T)^2)$. A more accurate transition model should therefore require fewer
 248 preference samples to achieve a given suboptimality, which we again confirm empirically.

249 7 Experimental results

250 We demonstrate the effectiveness of preference elicitation for offline reinforcement learning in prac-
 251 tice and compare the different sampling strategies introduced in Sections 5 and 6: OPRL with uni-
 252 form and uncertainty-sampling, and Sim-OPRL.

253 **Baselines.** For comparison, we also propose a practical implementation of Preference-based Opti-
 254 mistic Planning (PbOP), an uncertainty-based preference elicitation approach over trajectory rollouts
 255 in the *true environment* [Chen et al., 2022]. Finally, we report the performance of $\pi_{offline}^*$ and π^*
 256 as upper bounds for the performance of our algorithm: the former is trained in the learned transition
 257 model with access to the true reward, and the latter has full knowledge of both transition and reward
 258 function. We refer the reader to Appendix B for implementation details.

259 **Star MDP.** First, consider the tabular MDP in Figure 1a (we defer transition and reward details to
 260 Appendix C). Preferences collected over offline trajectories learn slowly about the negative reward in
 261 the bottom state, as it is always included in the comparison. Instead, simulated rollouts can directly
 262 query the optimal path. We thus find in Figure 1 that our preference elicitation strategy based on
 263 simulated rollouts achieves better returns than OPRL approaches, with fewer preference queries.

264 This example also illustrates the importance of pessimism with respect to the transition model. Even
 265 with access to true rewards, $\pi_{offline}^*$ avoids the out-of-distribution state, as it is unclear how to reach
 266 it. Thus, in Figure 1c, performance drops if pessimism is not applied to the output policy (purple
 267 lines). This confirms theoretical insights from Zhu et al. [2023], who demonstrate the importance
 268 of pessimism in offline preference-based RL. Pessimism is also crucial in simulated rollouts, to
 269 avoid wasting preference budget on regions of low confidence – as value estimates are in any case
 270 inaccurate. This is reflected in the lower efficiency of rollouts without pessimism over \hat{T} in Figure 1c
 271 (brown line). We also note the importance of optimism against reward uncertainty, both in OPRL in
 272 Figure 1b and in our model-based rollouts in Figure 1c.

273 Finally, as an upper bound for the performance of our algorithm, we include baselines that have
 274 access to the environment in Figure 1b: the optimal policy π^* , as well as an algorithm querying
 275 feedback over real environment rollouts [Chen et al., 2022, PbOP]. Final environment returns are
 276 higher than with Sim-OPRL, as they do not suffer from the limited coverage of the transition model.
 277 As supported by our theoretical analysis, this result stresses the importance of having a high-quality
 278 transition model to make our method effective. We explore this in more detail in the following.

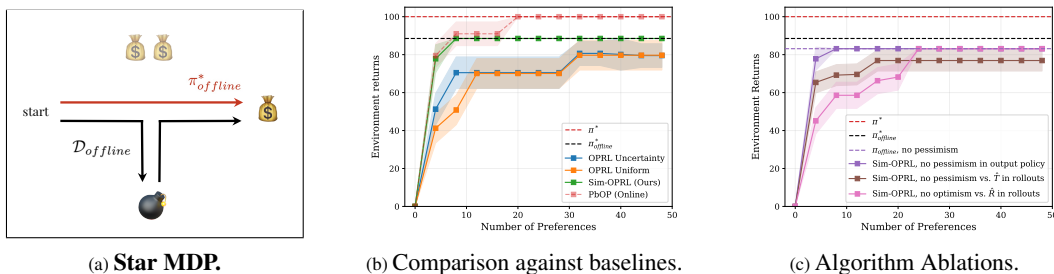


Figure 1: **Empirical results on the Star MDP.** Mean and 95% confidence interval over 20 experiments. Environment returns are normalized between 0 and 100. Only OPRL and Sim-OPRL are fully offline, all other methods have access to either environment interaction and/or to the true reward function.

Table 2: **Comparison of preference sample complexity N_p with different sampling methods**, to reach a suboptimality gap of $\epsilon = 20$ over normalized returns. Mean and 95% confidence interval over 20 experiments. The best-performing offline method is highlighted in bold.

Environment	OPRL Uniform	OPRL Uncertainty	Sim-OPRL (Ours)	PbOP (Online)
Star MDP (Figure 1a)	32 ± 4	30 ± 4	4 ± 2	4 ± 2
Gridworld	105 ± 11	66 ± 7	49 ± 7	32 ± 4
Sepsis Simulation	$18,856 \pm 427$	$2,246 \pm 143$	830 ± 88	261 ± 59

279 **Transition vs. Preference Model**
 280 **Quality.** Next, we study the trade-
 281 off between transition and preference
 282 model performance in our problem
 283 setting. In the low-data regime, eval-
 284 uation error due to the misspecifica-
 285 tion of the transition model is large.
 286 As dictated by our theoretical analys-
 287 is and as visualized in Figure 2a, this
 288 increases the number of preferences
 289 N_p required to achieve good final per-
 290 formance. Inversely, fewer prefer-
 291 ences are needed if the offline dataset
 292 is large and the transition model is ac-
 293 curate. We observe a similar trend for
 294 both Sim-OPRL and OPRL.

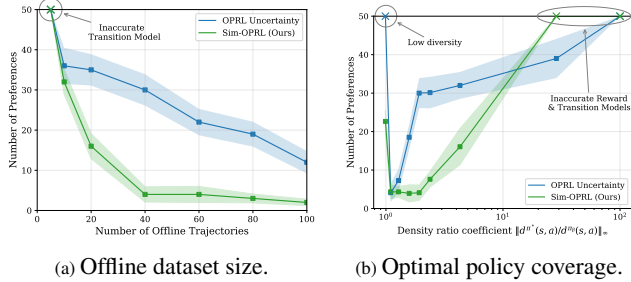


Figure 2: **Preference sample complexity N_p as function of the properties of the observational data**, to reach a suboptimality gap of $\epsilon = 20$ over normalized environment returns (Star MDP). Mean and 95% confidence intervals over 20 experiments. \times marks when the target suboptimality could not be achieved.

295 We also measure how coverage of the optimal policy affects performance. In Figure 2b, we vary the behavioral policy π_β underlying the offline data, ranging from optimal (density ratio 1) to suboptimal (large density ratio). We report the accuracy of transition and reward models in Appendix D. We observe that preference elicitation methods perform best when the data is close to optimal (except for a fully optimal, non-diverse dataset making reward learning from preferences challenging). More preference samples are required if the dataset has poor coverage of the optimal policy (large $C_T(\mathcal{F}_T, \pi^*)$), as transition and reward models become less accurate over the distribution of interest.

302 **Gridworld and Sepsis Simulation.** Finally, we validate our findings on more complex environments detailed in Appendix C: a gridworld experiment and a simulation of sepsis management in intensive care [Oberst and Sontag, 2019]. This example highlights another important advantage of Sim-OPRL over OPRL. In a sensitive setting such as healthcare where access is carefully controlled, it may be attractive to query experts about *synthetic* trajectories rather than real samples. Sample complexity results are given in Table 2, with similar conclusions: Sim-OPRL affords a higher preference sampling efficiency than OPRL baselines. For the sepsis environment, we note the number of preference samples needed to achieve our target suboptimality is large, likely due to the sparse nature of the reward function. In a real-world application, we could potentially warm-start the reward model by leveraging proxy rewards signals in the offline data [Yu et al., 2021].

312 8 Conclusion

313 Our work shows the potential of integrating human feedback within the framework of offline RL.
 314 We address the challenges of preference elicitation in a fully offline setup by exploring two key
 315 methods: sampling from the offline dataset [Shin et al., 2022, OPRL] and generating model rollouts
 316 (Sim-OPRL). By employing a pessimistic approach to handle out-of-distribution data and an opti-
 317 mistic strategy to acquire informative preferences, Sim-OPRL balances the need for robustness and
 318 informativeness in learning an optimal policy. We provide theoretical guarantees on the sample com-
 319 plexity of both approaches, demonstrating that performance depends on how well the offline data
 320 covers the optimal policy. Empirical evaluations in various environments confirm the effectiveness
 321 of our algorithm, as Sim-OPRL consistently outperforms baselines across different environments.

322 Overall, our approach not only advances the state-of-the-art in offline preference-based RL but also
 323 takes a significant step toward improving the practical utility of offline RL. This opens up new av-
 324 enues for real-world applications of RL in healthcare, robotics, and manufacturing, where interaction
 325 with the environment is challenging but domain experts can be queried for feedback.

326 References

- 327 Ralph Allan Bradley and Milton E Terry. Rank analysis of incomplete block designs: I. the method
328 of paired comparisons. *Biometrika*, 39(3/4):324–345, 1952.
- 329 Daniel Brown, Russell Coleman, Ravi Srinivasan, and Scott Niekum. Safe imitation learning
330 via fast Bayesian reward inference from preferences. In Hal Daumé III and Aarti Singh, ed-
331 itors, *Proceedings of the 37th International Conference on Machine Learning*, volume 119 of
332 *Proceedings of Machine Learning Research*, pages 1165–1177. PMLR, 13–18 Jul 2020. URL
333 <https://proceedings.mlr.press/v119/brown20a.html>.
- 334 Jonathan Chang, Masatoshi Uehara, Dhruv Sreenivas, Rahul Kidambi, and Wen Sun. Mitigating
335 covariate shift in imitation learning via offline data with partial coverage. *Advances in Neural*
336 *Information Processing Systems*, 34:965–979, 2021.
- 337 Xiaoyu Chen, Han Zhong, Zhuoran Yang, Zhaoran Wang, and Liwei Wang. Human-in-the-loop:
338 Provably efficient preference-based reinforcement learning with general function approximation.
339 In *International Conference on Machine Learning*, pages 3773–3793. PMLR, 2022.
- 340 Paul F Christiano, Jan Leike, Tom Brown, Miljan Martic, Shane Legg, and Dario Amodei. Deep
341 reinforcement learning from human preferences. *Advances in neural information processing sys-*
342 *tems*, 30, 2017.
- 343 Jonas Degraeve, Federico Felici, Jonas Buchli, Michael Neunert, Brendan Tracey, Francesco
344 Carpanese, Timo Ewalds, Roland Hafner, Abbas Abdolmaleki, Diego de Las Casas, et al. Mag-
345 netic control of tokamak plasmas through deep reinforcement learning. *Nature*, 602(7897):414–
346 419, 2022.
- 347 Steven Diamond and Stephen Boyd. Cvxpy: A python-embedded modeling language for convex
348 optimization. *Journal of Machine Learning Research*, 17(83):1–5, 2016.
- 349 Sara A Geer. *Empirical Processes in M-estimation*, volume 6. Cambridge university press, 2000.
- 350 Ying Jin, Zhuoran Yang, and Zhaoran Wang. Is pessimism provably efficient for offline rl? In
351 *International Conference on Machine Learning*, pages 5084–5096. PMLR, 2021.
- 352 Rahul Kidambi, Aravind Rajeswaran, Praneeth Netrapalli, and Thorsten Joachims. Morel: Model-
353 based offline reinforcement learning. *Advances in neural information processing systems*, 33:
354 21810–21823, 2020.
- 355 Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint*
356 *arXiv:1412.6980*, 2014.
- 357 Balaji Lakshminarayanan, Alexander Pritzel, and Charles Blundell. Simple and scalable predictive
358 uncertainty estimation using deep ensembles. *Advances in neural information processing systems*,
359 30, 2017.
- 360 Sergey Levine, Aviral Kumar, George Tucker, and Justin Fu. Offline reinforcement learning: Tuto-
361 rial, review, and perspectives on open problems. *arXiv preprint arXiv:2005.01643*, 2020.
- 362 David Lindner, Matteo Turchetta, Sebastian Tschiatschek, Kamil Ciosek, and Andreas Krause. In-
363 formation directed reward learning for reinforcement learning. *Advances in Neural Information*
364 *Processing Systems*, 34:3850–3862, 2021.
- 365 Qinghua Liu, Alan Chung, Csaba Szepesvári, and Chi Jin. When is partially observable reinforce-
366 ment learning not scary? In *Conference on Learning Theory*, pages 5175–5220. PMLR, 2022.
- 367 Azalia Mirhoseini, Anna Goldie, Mustafa Yazgan, Joe Jiang, Ebrahim Songhori, Shen Wang,
368 Young-Joon Lee, Eric Johnson, Omkar Pathak, Sungmin Bae, et al. Chip placement with deep
369 reinforcement learning. *arXiv preprint arXiv:2004.10746*, 2020.
- 370 Michael Oberst and David Sontag. Counterfactual off-policy evaluation with gumbel-max structural
371 causal models. In *International Conference on Machine Learning*, pages 4881–4890. PMLR,
372 2019.

- 373 Antonin Raffin, Ashley Hill, Adam Gleave, Anssi Kanervisto, Maximilian Ernestus, and Noah Dor-
374 mann. Stable-baselines3: Reliable reinforcement learning implementations. *Journal of Machine*
375 *Learning Research*, 22(268):1–8, 2021. URL [http://jmlr.org/papers/v22/20-1364.](http://jmlr.org/papers/v22/20-1364.html)
376 [html](http://jmlr.org/papers/v22/20-1364.html).
- 377 Aniruddh Raghu, Matthieu Komorowski, Imran Ahmed, Leo Celi, Peter Szolovits, and Marzyeh
378 Ghassemi. Deep reinforcement learning for sepsis treatment. *arXiv preprint arXiv:1711.09602*,
379 2017.
- 380 Dorsa Sadigh, Anca D Dragan, Shankar Sastry, and Sanjit A Seshia. Active preference-based learn-
381 ing of reward functions. 2018.
- 382 Aadirupa Saha, Aldo Pacchiano, and Jonathan Lee. Dueling rl: Reinforcement learning with tra-
383 jectory preferences. In *International Conference on Artificial Intelligence and Statistics*, pages
384 6263–6289. PMLR, 2023.
- 385 John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy
386 optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.
- 387 Daniel Shin, Anca Dragan, and Daniel S Brown. Benchmarks and algorithms for offline preference-
388 based reward learning. *Transactions on Machine Learning Research*, 2022.
- 389 Richard S Sutton and Andrew G Barto. *Reinforcement learning: An introduction*. MIT press, 2018.
- 390 Shengpu Tang and Jenna Wiens. Model selection for offline reinforcement learning: Practical con-
391 siderations for healthcare settings. In *Machine Learning for Healthcare Conference*, pages 2–35.
392 PMLR, 2021.
- 393 Masatoshi Uehara and Wen Sun. Pessimistic model-based offline reinforcement learning under
394 partial coverage. In *International Conference on Learning Representations*, 2021.
- 395 Christian Wirth, Riad Akrou, Gerhard Neumann, and Johannes Fürnkranz. A survey of preference-
396 based reinforcement learning methods. *Journal of Machine Learning Research*, 18(136):1–46,
397 2017.
- 398 Chao Yu, Jiming Liu, Shamim Nemati, and Guosheng Yin. Reinforcement learning in healthcare:
399 A survey. *ACM Computing Surveys (CSUR)*, 55(1):1–36, 2021.
- 400 Tianhe Yu, Garrett Thomas, Lantao Yu, Stefano Ermon, James Y Zou, Sergey Levine, Chelsea Finn,
401 and Tengyu Ma. Mopo: Model-based offline policy optimization. *Advances in Neural Information*
402 *Processing Systems*, 33:14129–14142, 2020.
- 403 Wenhao Zhan, Masatoshi Uehara, Nathan Kallus, Jason D Lee, and Wen Sun. Provable offline
404 reinforcement learning with human feedback. In *ICML 2023 Workshop The Many Facets of*
405 *Preference-Based Learning*, 2023a.
- 406 Wenhao Zhan, Masatoshi Uehara, Wen Sun, and Jason D Lee. How to query human feedback
407 efficiently in rl? In *ICML 2023 Workshop The Many Facets of Preference-Based Learning*,
408 2023b.
- 409 Banghua Zhu, Michael Jordan, and Jiantao Jiao. Principled reinforcement learning with human feed-
410 back from pairwise or k-wise comparisons. In *International Conference on Machine Learning*,
411 pages 43037–43067. PMLR, 2023.

412 **A Theoretical Details**

413 This appendix provides proofs for the presented theorems and lemmas. In subsection A.1, we pro-
 414 vide details on how we define the maximum likelihood estimators and confidence intervals of the
 415 preference and transition models. In subsection A.2 we provide the proof that our uncertainty-
 416 penalized objective in Equation (2) lower bounds the true value function and thus forms a valid
 417 pessimistic framework. In Appendix A.3, we show that the suboptimality of our offline preference
 418 elicitation framework is not lower-bounded by the performance of the optimal offline policy. In
 419 Appendix A.4, we provide our proof of theorem 5.1, analyzing the suboptimality of preferences
 420 sampled from an offline dataset. Finally, in Appendix A.5, we prove Theorem 6.1, which analyzes
 421 the suboptimality of preference sampling over simulated rollouts.

422 **A.1 Maximum Likelihood and Confidence Intervals**

423 Let \mathcal{F}_g denote a function class over $\mathcal{X} \rightarrow \Delta_{\mathcal{Y}}$, where \mathcal{X}, \mathcal{Y} are measurable sets, and $g \in \mathcal{F}_g$ denotes
 424 a function to be estimated.

425 Let \hat{g} denote the maximum likelihood estimator (MLE) of g based on a dataset $\mathcal{D} = \{(x^n, y^n)\}_{n=1}^N$:
 426 $\hat{g} = \operatorname{argmax}_{\tilde{g} \in \mathcal{F}_g} \mathbb{E}_{(x,y) \sim \mathcal{D}} \log(\tilde{g}(y|x))$. We construct the confidence set around the MLE as follows:

$$\mathcal{G} = \{\tilde{g} \in \mathcal{F}_g \mid \mathbb{E}_{(x,y) \sim \mathcal{D}} \left[\log \frac{\hat{g}(y|x)}{\tilde{g}(y|x)} \right] \leq \beta\}$$

Lemma A.1 (MLE Guarantee, Lemma 1 in Zhan et al. [2023a]). *Let $\delta \in (0, 1]$ and define the event \mathcal{E} that $g \in \mathcal{G}$. If*

$$\beta = \frac{c_{MLE} \log(\mathcal{N}_{\mathcal{F}_g}(1/N)/\delta)}{N},$$

427 *where $c_{MLE} > 0$ is a universal constant, then $P(\mathcal{E}) \geq 1 - \delta/2$.*

428 *Proof.* The proof follows that of Lemma 1 in Zhan et al. [2023a] and uses Cramér-Chernoff’s
 429 method.

430 Let $\bar{\mathcal{B}}$ be a $1/N$ -bracket of \mathcal{F}_g with $|\bar{\mathcal{B}}| = \mathcal{N}_{\mathcal{F}_g}(1/N)$. Denote the set of all right brackets in $\bar{\mathcal{B}}$ by
 431 $\tilde{\mathcal{B}} = \{b : \exists b' \text{ s.t. } [b', b] \in \bar{\mathcal{B}}\}$. For $b \in \tilde{\mathcal{B}}$, we have:

$$\begin{aligned} \mathbb{E} \left[\exp \left(\sum_{n=1}^N \log \frac{b(y^n|x^n)}{g(y^n|x^n)} \right) \right] &= \prod_{n=1}^N \mathbb{E} \left[\exp \left(\log \frac{b(y^n|x^n)}{g(y^n|x^n)} \right) \right] \\ &= \prod_{n=1}^N \mathbb{E} \left[\frac{b(y^n|x^n)}{g(y^n|x^n)} \right] \\ &= \prod_{n=1}^N \mathbb{E} \left[\sum_y b(y|x^n) \right] \\ &\leq (1 + 1/N)^N \leq e. \end{aligned}$$

432 as samples in \mathcal{D} as i.i.d. We use the Tower property in the third step and the fact that b is a $1/N$ -
 433 bracket for \mathcal{F}_g in the fourth: there exists $g' \in \mathcal{F}_g$ such that $\|g(\cdot|x) - b(\cdot|x)\|_1 \leq 1/N$ and thus
 434 $\|b(\cdot|x)\|_1 \leq 1 + 1/N$, for all $x \in \mathcal{X}$.

435 Then by Markov’s inequality, for any $\delta \in (0, 1]$, we have:

$$\begin{aligned} P \left(\sum_{n=1}^N \log \frac{b(y^n|x^n)}{g(y^n|x^n)} > \log(1/\delta) \right) &\leq \mathbb{E} \left[\exp \left(\sum_{n=1}^N \log \frac{b(y^n|x^n)}{g(y^n|x^n)} \right) \right] \cdot \exp(-\log(1/\delta)) \\ &\leq e\delta. \end{aligned}$$

436 By union bound, we have for all $b \in \tilde{\mathcal{B}}$,

$$P \left(\sum_{n=1}^N \log \frac{b(y^n|x^n)}{g(y^n|x^n)} > c_{MLE} \log(\mathcal{N}_{\mathcal{F}_g}(1/N)/\delta) \right) \leq \delta/2,$$

437 where $c_{MLE} > 0$ is a universal constant.

438 Finally, for all $\tilde{g} \in \mathcal{F}_g$, there exists $b \in \tilde{\mathcal{B}}$ such that $g(\cdot|x) \leq \tilde{g}(\cdot|x)$ for all $x \in \mathcal{X}$. As a result, for
 439 all $\tilde{g} \in \mathcal{F}_g$, we have:

$$P\left(\sum_{n=1}^N \log \frac{\tilde{g}(y^n|x^n)}{g(y^n|x^n)} > c_{MLE} \log(\mathcal{N}_{\mathcal{F}_g}(1/N)/\delta)\right) \leq \delta/2.$$

440

□

441 Under this event \mathcal{E} , we have $g \in \mathcal{G}$ with probability $1 - \delta/2$. A confidence interval constructed via
 442 loglikelihood also incurs a bound on the total variation (TV) distance between g and $\tilde{g} \in \mathcal{G}$:

443 **Lemma A.2** (TV-distance to MLE). *Under the event \mathcal{E} , we have, with probability $1 - \delta$, for all*
 444 $\tilde{g} \in \mathcal{G}$:

$$\mathbb{E}_{x \sim \mathcal{D}} [\|g(\cdot|x) - \tilde{g}(\cdot|x)\|_1^2] \leq \frac{c \log(\mathcal{N}_{\mathcal{F}_g}(1/N)/\delta)}{N}, \quad (4)$$

445 where $c > 0$ is a universal constant.

446 *Proof.* The proof follows that of Liu et al. [2022], Proposition 14. □

447 This guarantees that the true reward function is within an interval around the MLE estimate with
 448 high probability.

449 We apply these lemmas to our MLE estimates of transition and reward functions in Algorithm 1 to
 450 obtain the following guarantees.

451 Let \mathcal{E}_R denote the event $R \in \mathcal{R}$ and \mathcal{E}_T denote the event $T \in \mathcal{T}$, \mathcal{R} and \mathcal{T} denote the respective
 452 confidence sets around the MLE. By Lemma A.1, events \mathcal{E}_R and \mathcal{E}_T have probability $1 - \delta/2$ if we
 453 choose $\beta_R = c'_R \log(\mathcal{N}_{\mathcal{F}_R}(1/N_p)/\delta)/N_p$ and $\beta_T = c'_T \log(\mathcal{H}\mathcal{N}_{\mathcal{F}_T}(1/N_o)/\delta)/N_o$, where c'_R, c'_T
 454 are universal constants.

455 A.2 Model-based Pessimism and Uncertainty Penalties

456 **Lemma A.3** (Telescoping Lemma). *For any reward model $R \in \mathcal{F}_R$, and any two transition models*
 457 $T, \hat{T} \in \mathcal{F}_T$:

$$V_{T,R}^\pi - V_{\hat{T},R}^\pi \leq \mathbb{E}_{\tau \sim d_T^\pi} \left[\sum_{s_j, a_j \in \tau} \|T(\cdot|s_j, a_j) - \hat{T}(\cdot|s_j, a_j)\|_1 \right] \cdot R_{max}$$

458 *Proof.* The proof follows that of Lemma 4.1 in Yu et al. [2020] or Lemma 4 in Zhan et al. [2023a].

459 Let W_j be the expected return under policy π , with transition model \hat{T} for the first j steps, then
 460 transition model T for the rest of the episode. We have:

$$V_{T,R}^\pi - V_{\hat{T},R}^\pi = \sum_{j=0}^{H-1} W_j - W_{j+1}.$$

461 Now,

$$\begin{aligned} W_j &= R_j + \mathbb{E}_{s_j, a_j \sim \pi, \hat{T}} [\mathbb{E}_{s_{j+1} \sim T(\cdot|s_j, a_j)} [V_{T,R}^\pi(s_{j+1})]] \\ W_{j+1} &= R_j + \mathbb{E}_{s_j, a_j \sim \pi, \hat{T}} [\mathbb{E}_{s_{j+1} \sim \hat{T}(s_j, a_j)} [V_{T,R}^\pi(s_{j+1})]] \end{aligned}$$

462 where R_j is the expected return of the first j steps taken in \hat{T} . Therefore,

$$\begin{aligned} W_j - W_{j+1} &= \mathbb{E}_{s_j, a_j \sim \pi, \hat{T}} [\mathbb{E}_{s_{j+1} \sim T(\cdot|s_j, a_j)} [V_{T,R}^\pi(s_{j+1})] - \mathbb{E}_{s_{j+1} \sim \hat{T}(s_j, a_j)} [V_{T,R}^\pi(s_{j+1})]] \\ &\leq \mathbb{E}_{s_j, a_j \sim \pi, \hat{T}} [\|T(\cdot|s_j, a_j) - \hat{T}(\cdot|s_j, a_j)\|_1 \cdot R_{max}] \end{aligned}$$

463 under the boundedness assumption for R . Finally, we have:

$$\begin{aligned}
V_{T,R}^\pi - V_{\hat{T},R}^\pi &= \sum_{j=0}^{H-1} W_j - W_{j+1} \\
&= \sum_{j=0}^{H-1} \mathbb{E}_{s_j, a_j \sim \pi, \hat{T}} \left[\mathbb{E}_{s_{j+1} \sim T(\cdot | s_j, a_j)} [V_{T,R}^\pi(s_{j+1})] - \mathbb{E}_{s_{j+1} \sim \hat{T}(s_j, a_j)} [V_{T,R}^\pi(s_{j+1})] \right] \\
&\leq \sum_{j=0}^{H-1} \mathbb{E}_{s_j, a_j \sim \pi, \hat{T}} \left[\|T(\cdot | s_j, a_j) - \hat{T}(\cdot | s_j, a_j)\|_1 \cdot R_{max} \right] \\
&= \mathbb{E}_{\tau \sim d_T^\pi} \left[\sum_{s_j, a_j \in \tau} \|T(\cdot | s_j, a_j) - \hat{T}(\cdot | s_j, a_j)\|_1 \cdot R_{max} \right]
\end{aligned}$$

464

□

Lemma A.4 (Pessimistic Transition Model). *Under event \mathcal{E}_T , for all $\pi \in \Pi$, $\tilde{R} \in \mathcal{F}_R$:*

$$V_{\tilde{T}, \tilde{R} - u_T}^\pi \leq V_{T, \tilde{R}}^\pi.$$

Proof.

$$\begin{aligned}
V_{T, \tilde{R}}^\pi &= V_{\tilde{T}, \tilde{R}}^\pi - (V_{\tilde{T}, \tilde{R}}^\pi - V_{T, \tilde{R}}^\pi) \\
&\geq \mathbb{E}_{\tau \sim d_T^\pi} [\tilde{R}(\tau)] - \mathbb{E}_{\tau \sim d_T^\pi} [u_T(\tau)] \\
&= \mathbb{E}_{\tau \sim d_T^\pi} [\tilde{R}(\tau) - u_T(\tau)]
\end{aligned}$$

465 where we have used the telescoping lemma (Lemma A.3), and where $u_T(\tau) = \sum_{(s,a) \in \tau} u_T(s, a) \geq$
466 $\sum_{(s,a) \in \tau} \|\hat{T}(\cdot | s, a) - T(\cdot | s, a)\|_1 \cdot R_{max}$ under event \mathcal{E}_T . □

Lemma A.5 (Pessimistic Reward Model). *Under event \mathcal{E}_R , for all $\pi \in \Pi$, $\tilde{T} \in \mathcal{F}_T$:*

$$V_{\tilde{T}, \hat{R} - u_R}^\pi \leq V_{\tilde{T}, R}^\pi.$$

Proof.

$$\begin{aligned}
V_{\tilde{T}, R}^\pi &= V_{\tilde{T}, \hat{R}}^\pi - (V_{\tilde{T}, \hat{R}}^\pi - V_{\tilde{T}, R}^\pi) \\
&= \mathbb{E}_{\tau \sim d_{\tilde{T}}^\pi} [\hat{R}(\tau)] - \mathbb{E}_{\tau \sim d_{\tilde{T}}^\pi} [\hat{R}(\tau) - R(\tau)] \\
&\geq \mathbb{E}_{\tau \sim d_{\tilde{T}}^\pi} [\hat{R}(\tau) - u_R(\tau)]
\end{aligned}$$

467 where we have used the fact that $|\hat{R}(\tau) - R(\tau)| \leq \sum_{s,a \in \tau} |\hat{R}(s, a) - R(s, a)| =$
468 $\sum_{(s,a) \in \tau} u_R(s, a) = u_R(\tau)$ under event \mathcal{E}_R . □

469 Combining the above two lemmas gives the following result:

Corollary A.1. *Under events \mathcal{E}_T and \mathcal{E}_R , for all $\pi \in \Pi$:*

$$V_{\tilde{T}, \hat{R} - u_T - u_R}^\pi \leq V_{\tilde{T}, R}^\pi.$$

470 This justifies the overall objective considered in our pessimistic policy optimization procedure in
471 Section 4.

472 A.3 Suboptimality lower bound: a counterexample

473 Let $\pi_{offline}^* = \operatorname{argmax}_{\pi \in \Pi} \min_{\tilde{T} \in \mathcal{T}} V_{\tilde{T}, R}^\pi$ denote the optimal offline policy, which has access to
474 the ground-truth reward function. In this section, we ask whether its suboptimality $\epsilon_T = V_{T, R}^{\pi^*} -$
475 $V_{T, \hat{R}}^{\pi_{offline}^*}$ is a lower bound for the suboptimality of our learned policy $\hat{\pi}^*$ after preference elicitation.

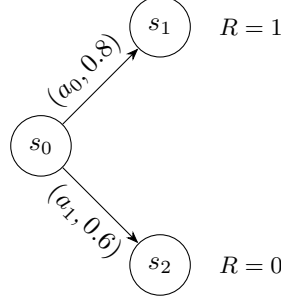


Figure 3: **Tabular MDP.** The environment starts in state s_0 and has horizon $H = 1$. Transition probabilities from state s_0 are given for the two binary actions a_0, a_1 (which send the agent to the other state with complementary probability).

476 **Counterexample.** Consider the MDP illustrated in Figure 3. Assume the following MLE estimate
477 and uncertainty function for both the transition and reward models:

$$\begin{aligned} \hat{T}(s_1|s_0, a_0) &= 0.5; & u_T(s_0, a_0) &= 0.4 \\ \hat{T}(s_1|s_0, a_1) &= 0.5; & u_T(s_0, a_1) &= 0.1 \\ \hat{r}(s_1) = \hat{r}(s_2) &= 0.5; & u_R(s_1) = u_R(s_2) &= 0.5 \end{aligned}$$

478 Assuming access to the learned transition model and the *true* reward function, we pessimistically
479 estimate the value of both actions:

$$\begin{aligned} V_{\hat{T}_{inf}, R}^{a_0} &= 0.1 \cdot 1 + 0.9 \cdot 0 = 0.1 \\ V_{\hat{T}_{inf}, R}^{a_1} &= 0.6 \cdot 0 + 0.4 \cdot 1 = 0.4 \end{aligned}$$

480 Thus, we have: $\pi_{offline}^*(s_0) = \operatorname{argmax}_a V_{\hat{T}_{inf}, R}^a = a_1$. The offline policy picks the suboptimal
481 action since the worst-case returns of this action are lower than those estimated for a_0 . Evaluating
482 this policy in the real environment, we get $\epsilon_T = V_{T, R}^{\pi^*} - V_{T, R}^{\pi_{offline}^*} = 0.6 \cdot 0 + 0.4 \cdot 1 = 0.4$.

483 We now estimate the optimal policy in the learned transition and reward model. Applying pessimism
484 with respect to both models, we get an equal estimated value of 0 for both actions a_0 and a_1 . If policy
485 optimization converges to $\hat{\pi}^* = a_0$, we reach the suboptimality $V_{T, R}^{\pi^*} - V_{T, R}^{\hat{\pi}^*} = 0.8 \cdot 1 + 0.2 \cdot 0 =$
486 $0.8 > \epsilon_T$.

487 This example demonstrates that ϵ_T is not a lower bound for the suboptimality of $\hat{\pi}^*$, as policy $\hat{\pi}^*$
488 can achieve lower suboptimality than $\pi_{offline}^*$ if errors in transition and reward model estimation
489 compensate each other.

490 A.4 Suboptimality of OPRL: Proof of Theorem 5.1

491 A.4.1 Suboptimality Decomposition

492 Recall that $\hat{T}_{inf}, \hat{R}_{inf} = \operatorname{argmin}_{\hat{T} \in \mathcal{T}, \hat{R} \in \mathcal{R}} V_{\hat{T}, \hat{R}}^\pi$ denote the pessimistic transition and reward models,
493 such that $\hat{\pi}^* = \operatorname{argmax}_{\pi \in \Pi} V_{\hat{T}_{inf}, \hat{R}_{inf}}^\pi$. We have:

$$\begin{aligned} V^{\pi^*} - V^{\hat{\pi}^*} &= V_{T, R}^{\pi^*} - V_{\hat{T}, \hat{R}}^{\hat{\pi}^*} \\ &= (V_{T, R}^{\pi^*} - V_{\hat{T}_{inf}, \hat{R}_{inf}}^{\pi^*}) - (V_{T, R}^{\hat{\pi}^*} - V_{\hat{T}_{inf}, \hat{R}_{inf}}^{\hat{\pi}^*}) \\ &\leq (V_{T, R}^{\pi^*} - V_{\hat{T}_{inf}, \hat{R}_{inf}}^{\pi^*}) - (V_{T, R}^{\hat{\pi}^*} - V_{\hat{T}_{inf}, \hat{R}_{inf}}^{\hat{\pi}^*}) \\ &\leq V_{T, R}^{\pi^*} - V_{\hat{T}_{inf}, \hat{R}_{inf}}^{\hat{\pi}^*}, \end{aligned} \tag{5}$$

494 where we have first used the optimality of $\hat{\pi}^*$ (stating that $V_{\hat{T}_{inf}, \hat{R}_{inf}}^{\hat{\pi}^*} \geq V_{\hat{T}_{inf}, \hat{R}_{inf}}^\pi$, for all π) and
495 then the pessimism principle (stating that $V_{\hat{T}_{inf}, \hat{R}_{inf}}^{\hat{\pi}^*} \leq V_{T, R}^{\hat{\pi}^*}$).

496 Finally, we decompose the last term above as follows:

$$V^{\pi^*} - V^{\tilde{\pi}^*} \leq \underbrace{(V_{T, \hat{R}_{inf}}^{\pi^*} - V_{\tilde{T}_{inf}, \hat{R}_{inf}}^{\pi^*})}_{\text{transition term}} + \underbrace{(V_{T, R}^{\pi^*} - V_{T, \hat{R}_{inf}}^{\pi^*})}_{\text{reward term}} \quad (6)$$

497 We further analyze each term in the following sections.

498 A.4.2 Analysis of the transition term

499 In this section, we now upper bound the transition term defined in Equation (6).

500 **Lemma A.6** (Lemma 4, Zhan et al. [2023a]). *Under the event \mathcal{E}_T , with probability $1 - \delta$, we have*
 501 *for all $\tilde{T} \in \mathcal{T}$, for all $\hat{R} \in \mathcal{G}_R$, for all π :*

$$\mathbb{E}_{d_T^\pi}[\tilde{R}(\tau)] - \mathbb{E}_{d_T^\pi}[\hat{R}(\tau)] \leq HR_{max} C_T(\mathcal{F}_T, \pi) \sqrt{\frac{c_T \log(H\mathcal{N}_{\mathcal{F}_T}(1/N_o)/\delta)}{N_o}},$$

502 where $c_T > 0$ is a constant.

503 *Proof.* From the telescoping lemma (Lemma A.3), we have:

$$\begin{aligned} V_{T, \hat{R}}^\pi - V_{\tilde{T}, \hat{R}}^\pi &\leq R_{max} \mathbb{E}_{\tau \sim d_T^\pi} \left[\sum_{s_j, a_j \in \tau} \|T(\cdot|s_j, a_j) - \tilde{T}(\cdot|s_j, a_j)\|_1 \right] \\ &\leq HR_{max} \mathbb{E}_{(s, a) \sim d_T^\pi} [\|T(\cdot|s, a) - \tilde{T}(\cdot|s, a)\|_1] \\ &\leq HR_{max} C_T(\mathcal{F}_T, \pi) \sqrt{\mathbb{E}_{(s, a) \sim D_{offline}} [\|T(\cdot|s, a) - \tilde{T}(\cdot|s, a)\|_1^2]} \end{aligned}$$

Under event \mathcal{E}_T , by Lemma A.2, we have, with probability $1 - \delta$, for all $\tilde{T} \in \mathcal{T}$:

$$\mathbb{E}_{(s, a) \sim D_{offline}} [\|T(\cdot|s, a) - \tilde{T}(\cdot|s, a)\|_1^2] \leq \frac{1}{N_o} c_T \log(H\mathcal{N}_{\mathcal{F}_T}(1/N_o)/\delta)$$

504 This concludes our proof.

505 □

506 A.4.3 Analysis of the reward term

507 Next, we upper bound the reward term defined in Equation (6).

508 As in Zhan et al. [2023a], we consider the following value function: $V_{T, R}^\pi = \mathbb{E}_{\tau \sim d_T^\pi} [R(\tau)] -$
 509 $\mathbb{E}_{\tau \sim d_{pref}} [R(\tau)]$, where d_{pref} is a fixed reference trajectory distribution. This baseline subtraction,
 510 which doesn't affect either the optimal policy or the analysis of the transition term, is needed as the
 511 approximated confidence set is based on the uncertainty in *preference* between two trajectories, not
 512 in the reward of a single one.

513 **Definition A.1** (Preference concentrability coefficient). *The concentrability coefficient w.r.t. reward*
 514 *classes \mathcal{F}_R , a target policy π^* and a reference trajectory distribution d_{pref} is defined as:*

$$C_R(\mathcal{F}_R, \pi^*) = \frac{\mathbb{E}_{\tau_1 \sim d_T^{\pi^*}, \tau_2 \sim d_{pref}} [u_{P_R}(\tau_1, \tau_2)]}{\mathbb{E}_{\tau_1, \tau_2 \sim D_{offline}} [u_{P_R}(\tau_1, \tau_2)]}$$

515 Note that, for the purpose of our analysis, our definition differs from that of Zhan et al. [2023a] who
 516 instead consider the max ratio of difference in rewards term: $|R(\tau_1) - R(\tau_2) - \hat{R}(\tau_1) + \hat{R}(\tau_2)|$ over
 517 the entire function class \mathcal{F}_R .

518 **Lemma A.7.** *Let trajectories for preference elicitation be sampled uniformly from the offline*
 519 *dataset. Under the event \mathcal{E}_R , with probability $1 - \delta$, we have for all $\tilde{T} \in \mathcal{G}_T$, for all $\hat{R} \in \mathcal{R}$,*
 520 *for all π :*

$$V_{T, R}^{\pi^*} - V_{T, \hat{R}_{inf}}^{\pi^*} \leq 2\kappa C_R(\mathcal{F}_T, \pi) \sqrt{\frac{c_R \log(\mathcal{N}_{\mathcal{F}_R}(1/N_p)/\delta)}{N_p}},$$

521 where $c_R > 0$ is a constant and $\kappa = \sup_{r \in [-R_{max}, R_{max}]} \frac{1}{\sigma'(r)}$ measures the degree of non-linearity
 522 of the sigmoid function.

Proof.

$$\begin{aligned}
 V_{T,R}^{\pi^*} - V_{T,\hat{R}_{inf}}^{\pi^*} &= \mathbb{E}_{\tau \sim d_T^{\pi^*}} [R(\tau)] - \mathbb{E}_{\tau \sim d_{pref}} [R(\tau)] - \mathbb{E}_{\tau \sim d_T^{\pi^*}} [\hat{R}_{inf}(\tau)] + \mathbb{E}_{\tau \sim d_{pref}} [\hat{R}_{inf}(\tau)] \\
 &= \mathbb{E}_{\tau_1 \sim d_T^{\pi^*}, \tau_2 \sim d_{pref}} [R(\tau_1) - R(\tau_2)] - (\hat{R}_{inf}(\tau_1) - \hat{R}_{inf}(\tau_2)) \\
 &\leq \kappa \mathbb{E}_{\tau_1 \sim d_T^{\pi^*}, \tau_2 \sim d_{pref}} [|P_R(\tau_1 \succ \tau_2) - P_{\hat{R}_{inf}}(\tau_1 \succ \tau_2)|] \\
 &\leq \kappa \mathbb{E}_{\tau_1 \sim d_T^{\pi^*}, \tau_2 \sim d_{pref}} [u_{P_R}(\tau_1, \tau_2)] \\
 &= \kappa C'_R(\mathcal{F}_T, \pi^*) \mathbb{E}_{\tau_1, \tau_2 \sim \mathcal{D}_{offline}} [u_{P_R}(\tau_1, \tau_2)] \tag{7}
 \end{aligned}$$

523 where $\kappa = \sup_{r \in [-R_{max}, R_{max}]} \frac{1}{\sigma'(r)}$ measures the degree of non-linearity of the sigmoid function.
 524 In the first inequality, we have applied the mean value theorem, under Assumption 3.2. In the second
 525 inequality, we have used the definition of uncertainty function u_{P_R} as we know $\hat{R}_{inf} \in \mathcal{R}$.

526 Now, under event \mathcal{E}_R , by Lemma A.2, we have, with probability $1 - \delta$ for all $\tilde{R} \in \mathcal{R}$:

$$\mathbb{E}_{(\tau_1, \tau_2) \sim \mathcal{D}_{pref}} [\|P_R(\tau_1 \succ \tau_2) - P_{\tilde{R}}(\tau_1 \succ \tau_2)\|_1^2] \leq \frac{c_R \log(\mathcal{N}_{\mathcal{F}_R}(1/N_p)/\delta)}{N_p}, \tag{8}$$

527 where $c_R > 0$ is a constant. This implies the following upper bound for the preference uncertainty
 528 function:

$$\mathbb{E}_{(\tau_1, \tau_2) \sim \mathcal{D}_{pref}} [u_{P_R}(\tau_1, \tau_2)] \leq 2 \sqrt{\frac{c_R \log(\mathcal{N}_{\mathcal{F}_g}(1/N_p)/\delta)}{N_p}} \tag{9}$$

529 Under uniform sampling, the distribution of preferences in \mathcal{D}_{pref} is that of the offline dataset:

$$\mathbb{E}_{(\tau_1, \tau_2) \sim \mathcal{D}_{offline}} [u_{P_R}(\tau_1, \tau_2)] = \mathbb{E}_{(\tau_1, \tau_2) \sim \mathcal{D}_{pref}} [u_{P_R}(\tau_1, \tau_2)]$$

530 Thus,

$$V_{T,R}^{\pi^*} - V_{T,\hat{R}_{inf}}^{\pi^*} \leq 2\kappa C'_R(\mathcal{F}_T, \pi) \sqrt{\frac{c_R \log(\mathcal{N}_{\mathcal{F}_R}(1/N_p)/\delta)}{N_p}}.$$

531

□

532 **Lemma A.8.** *Let trajectories for preference elicitation be sampled through uncertainty sampling*
 533 *from the offline dataset. Under the event \mathcal{E}_R , with probability $1 - \delta$, we have for all $\tilde{T} \in \mathcal{G}_T$, for all*
 534 *$\tilde{R} \in \mathcal{R}$, for all π :*

$$V_{\tilde{T},R}^{\pi^*} - V_{\tilde{T},\hat{R}_{inf}}^{\pi^*} \leq 2\alpha\kappa C'_R(\mathcal{F}_T, \pi) \sqrt{\frac{c_R \log(\mathcal{N}_{\mathcal{F}_R}(1/N_p)/\delta)}{N_p}},$$

535 where $c_R > 0$ is a constant and $\alpha \leq 1$.

Proof. The proof follows closely that of Lemma A.7. We introduce the preference concentrability coefficient defined for a general preference dataset:

$$C'_R(\mathcal{F}_R, \pi^*) = \frac{\mathbb{E}_{\tau_1 \sim d_T^{\pi^*}, \tau_2 \sim d_{pref}} [u_{P_R}(\tau_1, \tau_2)]}{\mathbb{E}_{\tau_1, \tau_2 \sim \mathcal{D}_{pref}} [u_{P_R}(\tau_1, \tau_2)]}$$

536 We start from Equation (7):

$$\begin{aligned}
 V_{T,R}^{\pi^*} - V_{T,\hat{R}_{inf}}^{\pi^*} &\leq \kappa \mathbb{E}_{\tau_1 \sim d_T^{\pi^*}, \tau_2 \sim d_{pref}} [u_{P_R}(\tau_1, \tau_2)] \\
 &= \kappa C'_R(\mathcal{F}_T, \pi^*) \mathbb{E}_{\tau_1, \tau_2 \sim \mathcal{D}_{pref}} [u_{P_R}(\tau_1, \tau_2)] \\
 &\leq 2\kappa C'_R(\mathcal{F}_T, \pi^*) \sqrt{\frac{c_R \log(\mathcal{N}_{\mathcal{F}_g}(1/N_p)/\delta)}{N_p}}
 \end{aligned}$$

537 where we have used Equation (9).

Now consider the dataset of uncertainty-sampled preferences \mathcal{D}_{pref} . By definition, we have:

$$\mathbb{E}_{\tau_1, \tau_2 \sim \mathcal{D}_{pref}} [u_{P_R}(\tau_1, \tau_2)] \geq \mathbb{E}_{\tau_1, \tau_2 \sim \mathcal{D}_{offline}} [u_{P_R}(\tau_1, \tau_2)]$$

538 Thus, we have: $C'_R(\mathcal{F}_T, \pi^*) \leq C_R(\mathcal{F}_T, \pi^*)$. In other words, we can write: $C'_R(\mathcal{F}_T, \pi^*) =$
 539 $\alpha C_R(\mathcal{F}_T, \pi^*)$, where $\alpha \leq 1$. This concludes our proof. \square

540 We now conclude the proof of Theorem 5.1 under events \mathcal{E}_R and \mathcal{E}_T .

From Lemma A.6, we upper bound the transition term:

$$V_{T, \hat{R}_{inf}}^{\pi^*} - V_{\hat{T}_{inf}, \hat{R}_{inf}}^{\pi^*} \leq HR_{max} C_T(\mathcal{F}_T, \pi^*) \sqrt{\frac{c_T \log(HN_{\mathcal{F}_T}(1/N_o)/\delta)}{N_o}}$$

From Lemmas A.7 and A.8, we upper bound the reward term:

$$V_{T, R}^{\pi^*} - V_{T, \hat{R}_{inf}}^{\pi^*} \leq 2\alpha\kappa C_R(\mathcal{F}_T, \pi^*) \sqrt{\frac{c_R \log(N_{\mathcal{F}_R}(1/N_p)/\delta)}{N_p}},$$

541 where $\alpha = 1$ for uniform sampling or $\alpha \leq 1$ for uncertainty sampling.

542 Combining with Equation (6), we obtain Theorem 5.1.

543 A.5 Suboptimality of Sim-OPRL: Proof of Theorem 6.1

544 A.5.1 Suboptimality Decomposition

545 We decompose the suboptimality slightly differently to Equation (5), introducing the optimal
 546 offline policy (optimal in the pessimistic model under the *true* reward function): $\pi_{offline}^* =$
 547 $\operatorname{argmax}_{\pi \in \Pi} V_{\hat{T}_{inf}, R}^{\pi}$.

$$\begin{aligned} V^{\pi^*} - V^{\hat{\pi}^*} &= V_{T, R}^{\pi^*} - V_{\hat{T}, R}^{\hat{\pi}^*} \\ &= (V_{T, R}^{\pi^*} - V_{\hat{T}_{inf}, \hat{R}_{inf}}^{\pi_{offline}^*}) - (V_{\hat{T}, R}^{\hat{\pi}^*} - V_{\hat{T}_{inf}, \hat{R}_{inf}}^{\pi_{offline}^*}) \\ &\leq (V_{T, R}^{\pi^*} - V_{\hat{T}_{inf}, \hat{R}_{inf}}^{\pi_{offline}^*}) - (V_{\hat{T}, R}^{\hat{\pi}^*} - V_{\hat{T}_{inf}, \hat{R}_{inf}}^{\hat{\pi}^*}) \\ &\leq V_{T, R}^{\pi^*} - V_{\hat{T}_{inf}, \hat{R}_{inf}}^{\pi_{offline}^*} \\ &= (V_{T, R}^{\pi^*} - V_{\hat{T}_{inf}, R}^{\pi^*}) + (V_{\hat{T}_{inf}, R}^{\pi^*} - V_{\hat{T}_{inf}, \hat{R}_{inf}}^{\pi_{offline}^*}) \\ &\leq \underbrace{(V_{T, R}^{\pi^*} - V_{\hat{T}_{inf}, R}^{\pi^*})}_{\text{transition term}} + \underbrace{(V_{\hat{T}_{inf}, R}^{\pi^*} - V_{\hat{T}_{inf}, \hat{R}_{inf}}^{\pi_{offline}^*})}_{\text{reward term}} \end{aligned} \quad (10)$$

548 where we have followed the same analysis as in Appendix A.4.1 and used the optimality of $\pi_{offline}^*$
 549 in the last inequality.

550 The analysis of the transition term is identical to the above (Appendix A.4.2). We analyze the reward
 551 term next.

552 A.5.2 Analysis of the reward term

553 **Lemma A.9** (Optimal Offline Policy In Set). *Let $\Pi_{offline}$ denote the following set of near-optimal*
 554 *pessimistic policies, under the pessimistic transition model \hat{T}_{inf} and the reward confidence set \mathcal{R} :*

$$\Pi_{offline} = \{\pi \mid \pi = \operatorname{argmax}_{\pi \in \Pi} \mathbb{E}_{\tau \sim d_{\hat{T}_{inf}}^{\pi}} [\tilde{R}(\tau)] \forall \tilde{R} \in \mathcal{R}\}$$

555 Under event \mathcal{E}_R , we have $\pi_{offline}^* \in \Pi_{offline}$.

556 *Proof.* Recall the definition of $\pi_{offline}^*$: $\pi_{offline}^* = \operatorname{argmax}_{\pi \in \Pi} V_{\hat{T}_{inf}, R}^\pi$. Note that there is no need
 557 to consider the preference baseline term in V^π when building $\Pi_{offline}$ since it is independent of the
 558 policy. Under event \mathcal{E}_R , we have $R \in \mathcal{R}$. Thus, $\pi_{offline}^* \in \Pi_{offline}$. \square

559 **Lemma A.10.** *Under event \mathcal{E}_R , we have, with probability $1 - \delta$:*

$$V_{\hat{T}_{inf}, R}^{\pi_{offline}^*} - V_{\hat{T}_{inf}, \hat{R}_{inf}}^{\pi_{offline}^*} \leq 2\sqrt{\kappa^2 c_R / N_p \log(\mathcal{N}_{\mathcal{F}_R}(1/N_p)/\delta)}$$

Proof.

$$\begin{aligned} & V_{\hat{T}_{inf}, R}^{\pi_{offline}^*} - V_{\hat{T}_{inf}, \hat{R}_{inf}}^{\pi_{offline}^*} \\ &= (V_{\hat{T}_{inf}, R}^{\pi_{offline}^*} - V_{\hat{T}_{inf}, \hat{R}}^{\pi_{offline}^*}) + (V_{\hat{T}_{inf}, \hat{R}}^{\pi_{offline}^*} - V_{\hat{T}_{inf}, \hat{R}_{inf}}^{\pi_{offline}^*}) \\ &= \mathbb{E}_{\tau \sim d_{\hat{T}_{inf}}^{\pi_{offline}^*}} [R(\tau)] - \mathbb{E}_{\tau \sim d_{pref}} [R(\tau)] - \mathbb{E}_{\tau \sim d_{\hat{T}_{inf}}^{\pi_{offline}^*}} [\hat{R}_{inf}(\tau)] + \mathbb{E}_{\tau \sim d_{pref}} [\hat{R}_{inf}(\tau)] \\ &= \mathbb{E}_{\tau_1 \sim d_{\hat{T}_{inf}}^{\pi_{offline}^*}, \tau_2 \sim d_{pref}} [R(\tau_1) - R(\tau_2)] - \mathbb{E}_{\tau_1 \sim d_{\hat{T}_{inf}}^{\pi_{offline}^*}, \tau_2 \sim d_{pref}} [\hat{R}_{inf}(\tau_1) - \hat{R}_{inf}(\tau_2)] \\ &\leq \kappa \mathbb{E}_{\tau_1 \sim d_{\hat{T}_{inf}}^{\pi_{offline}^*}, \tau_2 \sim d_{pref}} [P_R(\tau_1 \succ \tau_2) - P_{\hat{R}_{inf}}(\tau_1 \succ \tau_2)], \end{aligned}$$

560 where $\kappa = \sup_{r \in [-R_{max}, R_{max}]} \frac{1}{\sigma'(r)}$ measures the degree of non-linearity of the sigmoid function.

561 We have applied the mean value theorem, under Assumption 3.2.

562 As $R_{inf} \in \mathcal{R}$, we have: $P_R(\tau_1 \succ \tau_2) - P_{\hat{R}_{inf}}(\tau_1 \succ \tau_2) \leq u_{P_R}(\tau_1, \tau_2)$.

563 Let d_{pref} correspond to the distribution of the preference data, which consists of rollouts from
 564 exploratory policies within the learned environment model: $d_{pref} = d_{\hat{T}_{inf}}^{\pi_1} / 2 + d_{\hat{T}_{inf}}^{\pi_2} / 2$. Recall that

565 the near-optimal policy set $\Pi_{offline}$ includes policy $\pi_{offline}^*$ (Lemma A.9) and that π_1, π_2 are the
 566 two more exploratory policies within this set:

$$\mathbb{E}_{\tau_1 \sim d_{\hat{T}_{inf}}^{\pi_{offline}^*}, \tau_2 \sim d_{pref}} [u_{P_R}(\tau_1, \tau_2)] \leq \max_{\pi_1, \pi_2 \in \Pi_{offline}} \mathbb{E}_{\tau_1 \sim d_{\hat{T}_{inf}}^{\pi_1}, \tau_2 \sim d_{\hat{T}_{inf}}^{\pi_2}} [u_{P_R}(\tau_1, \tau_2)].$$

567 Now, under event \mathcal{E}_R , by Lemma A.2, we have, with probability $1 - \delta$ for all $\tilde{R} \in \mathcal{R}$:

$$\mathbb{E}_{(\tau_1, \tau_2) \sim \mathcal{D}_{pref}} [\|P_R(\tau_1 \succ \tau_2) - P_{\tilde{R}}(\tau_1 \succ \tau_2)\|_1^2] \leq \frac{c_R \log(\mathcal{N}_{\mathcal{F}_R}(1/N_p)/\delta)}{N_p},$$

568 where $c_R > 0$ is a constant. This implies the following upper bound for the preference uncertainty
 569 function:

$$\mathbb{E}_{(\tau_1, \tau_2) \sim \mathcal{D}_{pref}} [u_{P_R}(\tau_1, \tau_2)] \leq 2\sqrt{\frac{c_R \log(\mathcal{N}_{\mathcal{F}_g}(1/N_p)/\delta)}{N_p}}.$$

570 Thus, we obtain:

$$V_{\hat{T}_{inf}, R}^{\pi_{offline}^*} - V_{\hat{T}_{inf}, \hat{R}_{inf}}^{\pi_{offline}^*} \leq 2\kappa \sqrt{\frac{c_R \log(\mathcal{N}_{\mathcal{F}_g}(1/N_p)/\delta)}{N_p}}.$$

571 The resulting sample complexity of $\mathcal{O}(\frac{\kappa^2 d}{\epsilon^2})$ matches that of active preference learning within a
 572 known environment [Saha et al., 2023, Chen et al., 2022].

573 \square

574 We now conclude the proof of Theorem 6.1 under events \mathcal{E}_R and \mathcal{E}_T .

From Lemma A.6, we upper bound the transition term:

$$V_{T, R}^\pi - V_{\hat{T}_{inf}, R}^{\pi^*} \leq H R_{max} C_T(\mathcal{F}_T, \pi^*) \sqrt{\frac{c_T \log(H \mathcal{N}_{\mathcal{F}_T}(1/N_o)/\delta)}{N_o}}.$$

From Lemma A.10, we upper bound the reward term:

$$V_{\hat{T}_{inf,R}}^{\pi_{offline}^*} - V_{\hat{T}_{inf},\hat{R}_{inf}}^{\pi_{offline}^*} \leq 2\kappa \sqrt{\frac{c_R \log(\mathcal{N}_{\mathcal{F}_R}(1/N_p)/\delta)}{N_p}}.$$

575 Combining with Equation (10), we obtain Theorem 6.1.

576 B Implementation Details

577 We trained all models on two 64-core AMD processors or a single NVIDIA RTX2080Ti GPU. The
578 total wall-clock time for running all experiments presented in this paper amounted to less than 72
579 hours.

580 **Transition and Reward Function Training.** For all baselines, transition and reward models were
581 implemented as linear classifiers (for the Star MDP) or as two-layer perceptions with ReLU activa-
582 tion and hidden layer dimension 32 (Gridworld and Sepsis environments). Training was carried out
583 for two or one epochs for the transition and reward models respectively, with the Adam optimizer
584 [Kingma and Ba, 2014] and a learning rate of 10^{-3} .

585 We provide a more detailed practical algorithm for Sim-OPRL in Algorithm 3. For both our method
586 and baselines relying on uncertainty sets (OPRL and PbOP), we estimated uncertainty sets by train-
587 ing models initialized with different random seeds on different bootstraps of the data (sampling 90%
588 of the data with replacement). We consider ensembles of size $|\mathcal{T}| = |\mathcal{R}| = 5$ for both transition and
589 reward models. Hyperparameters λ_T, λ_R control the degree of pessimism in practice and could be
590 considered equivalent to adjusting margin parameters β_T, β_R in our conceptual algorithm proposed
591 in Section 4. Since the exact values prescribed by our theoretical analysis cannot be estimated, the
592 user must set these parameters themselves. Hyperparameter optimization in offline RL is a chal-
593 lenging problem [Levine et al., 2020]; for our experiments, we simply set $\lambda_T = 0.5, \lambda_R = 0.1$
594 (StarMDP, Gridworld) and $\lambda_T = \lambda_R = 1$ for the Sepsis environment.

Algorithm 3 Sim-OPRL: Practical Algorithm

Input: Observational trajectories dataset $\mathcal{D}_{offline}$. Hyperparameters λ_T, λ_R .

Output: $\hat{\pi}^*$

1: Train an ensemble \mathcal{T} of transition models via bootstrapping on the observational data $\mathcal{D}_{offline}$:

$$\hat{T}(\cdot|s, a) = \frac{1}{|\mathcal{T}|} \sum_{\tilde{T} \in \mathcal{T}} \tilde{T}(\cdot|s, a); \quad u_T(s, a) = \max_{T_1, T_2 \in \mathcal{T}} |T_1(\cdot|s, a) - T_2(\cdot|s, a)| \cdot R_{max}$$

2: $\mathcal{D}_{pref} \leftarrow \emptyset$.

3: **for** $k = 1, \dots, N_p$ **do**

4: Estimate optimal offline policy set:

$$\Pi_{offline} = \{\pi \mid \pi = \operatorname{argmax}_{\pi \in \Pi} \mathbb{E}_{(s,a) \sim d_T^\pi} [\tilde{R}(s, a) - \lambda_T u_T(s, a)] \forall \tilde{R} \in \mathcal{R}\}$$

5: Identify exploratory policies: $\pi_1, \pi_2 = \operatorname{argmax}_{\pi_1, \pi_2 \in \Pi_{offline}} \mathbb{E}_{\tau_1 \sim d_T^{\pi_1}, \tau_2 \sim d_T^{\pi_2}} [u_{P_R}(\tau_1, \tau_2)]$

6: Rollouts in model: $\tau_1 \sim d_T^{\pi_1}, \tau_2 \sim d_T^{\pi_2}$.

7: Collect preference label o for (τ_1, τ_2) .

8: $\mathcal{D}_{pref} \leftarrow \mathcal{D}_{pref} \cup \{(\tau_1, \tau_2, o)\}$.

9: Train an ensemble \mathcal{R} of reward models via bootstrapping of the preference data \mathcal{D}_{pref} :

$$\hat{R}(s, a) = \frac{1}{|\mathcal{R}|} \sum_{\tilde{R} \in \mathcal{R}} \tilde{R}(s, a); \quad u_R(s, a) = \max_{R_1, R_2 \in \mathcal{R}} |R_1(\cdot|s, a) - R_2(\cdot|s, a)|$$

10: **end for**

11: $\hat{\pi}^* \leftarrow \operatorname{argmax}_{\pi \in \Pi} \mathbb{E}_{(s,a) \sim d_T^\pi} [\hat{R}(s, a) - \lambda_R u_R(s, a) - \lambda_T u_T(s, a)]$

595 **Near-Optimal Policy Set and Exploratory Policies.** Both Sim-OPRL and PbOP require con-
 596 structing a set of near-optimal policies within a learned model of the environment. Note that the
 597 PbOP algorithm in Chen et al. [2022] proposes to construct the near-optimal policy set by consider-
 598 ing all policies that have a preference greater than $1/2$ over *all other policies* in Π , under a transition
 599 and preference uncertainty bonus. This is infeasible to estimate in practice; we modified the algo-
 600 rithm to allow for practical implementation. The motivation in building the set of plausibly optimal
 601 policies remains the same, but the theoretical guarantees may not hold.

602 We build $\Pi_{offline}$ by maintaining a policy model for all $\tilde{R} \in \mathcal{R}$, i.e., each element of the reward
 603 ensemble. Policy models are optimized to maximize returns under the transition model \hat{T} and the
 604 reward function $\tilde{R} - \lambda_T u_T$ (Sim-OPRL) or $\tilde{R} + \lambda_T u_T$ (PbOP). Next, the most exploratory poli-
 605 cies are identified by generating 10 rollouts of each of the candidate policies within the learned
 606 (SimOPRL) or true (PbOP) model. The trajectories (τ_1, τ_2) maximizing the preference uncertainty
 607 function $u_{P_{\tilde{R}}}(\tau_1, \tau_2)$ are used for preference feedback. In PbOP, the trajectories are then added to
 608 the trajectories buffer and the transition model is retrained for 20 (Star MDP, Gridworld) or 200
 609 steps (Sepsis).

610 **Preference Feedback Collection.** Preference labels are provided through the ground-truth reward
 611 function associated with every environment. As stated in Section 4, for computational efficiency, we
 612 sample preferences in batches of 4 (Star MDP, Gridworld) or 100 (Sepsis) to reduce the number of
 613 model updates needed.

614 **Policy Optimization.** Policy optimization stages, both in estimating optimal policy sets in Sim-
 615 OPRL and PbOP and in outputting final policies, are carried out exactly through linear programming
 616 for the Star MDP and Gridworld using `cvxopt` [Diamond and Boyd, 2016], based on code from
 617 Lindner et al. [2021], and using Proximal Policy Optimization [Schulman et al., 2017] implemented
 618 in `stable-baselines3` [Raffin et al., 2021] for the Sepsis environment. In the latter case, after
 619 every preference collection episode, reward and policy models were trained from the checkpoint of
 620 the previous iteration, for only 20 steps to minimize computation.

621 **Baselines and Ablations.** We implement both OPRL baselines within our model-based offline
 622 preference-based algorithm described in Section 4. Uncertainty sampling is taking the pair with
 623 maximum preference uncertainty over 45 pairs for every sample, to reduce the load of computing
 624 preference uncertainty over the entire trajectory buffer.

625 Our ablation study for Figure 1c is conducted as follows. For Sim-OPRL without pessimism in the
 626 output policy, we output the policy that maximizes the value function under the MLE estimate of
 627 the transition and reward function, \hat{T} and \hat{R} , after preference acquisition. For Sim-OPRL without
 628 pessimism in the simulated rollouts, we estimate the optimal policy set $\Pi_{offline}$ in the MLE esti-
 629 mate of the transition model instead of its pessimistic counterpart. Finally, for Sim-OPRL without
 630 optimism in the simulated rollouts, we generate rollouts from any two policies in $\Pi_{offline}$ instead
 631 of the most explorative ones.

632 C Environment Details

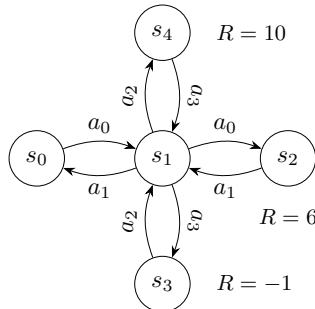


Figure 4: **Star MDP illustrated in Figure 1a.** Transition probabilities are 0.9 for all solid arrows. Omitted actions or complementary transitions keep the state unchanged.

633 **Star MDP.** We illustrate the transition dynamics underlying the Star MDP in Figure 4. Transition
 634 probabilities are 0.9 for all depicted solid arrows, and leave the state unchanged otherwise. Other
 635 actions also keep the state unchanged with probability 1. Episodes have length $H = 3$ and start
 636 from s_0 . Unless specified otherwise, the offline dataset $\mathcal{D}_{offline}$ consists of 40 trajectories which
 637 only cover states (s_0, s_1, s_3) and (s_3, s_1, s_2) .

Start		-1	-1
		-1	-1
	20		
			10

Figure 5: **Gridworld environment.** Rewards at every state are indicated if non-zero. Transition probabilities are 0.9. Thick lines indicate an obstacle, through which state transitions have probability zero.

638 **Gridworld.** We illustrate the gridworld environment in Figure 5. The environment consists of a
 639 4×4 grid with states associated with different rewards, including a negative-reward region in the
 640 top-right corner, a high-reward but unreachable state, and a moderate-reward state at the bottom
 641 right corner. Each episode starts in the top-left corner. Transition probabilities for each of the four
 642 actions (top, left, bottom, right) are 0.9 for the intended direction, and 0.1 for the others;
 643 and action stay remains in the current state with probability 1. Transitions beyond the grid limits or
 644 through obstacles have probability zero, with the remainder of the probability mass for each action
 645 being distributed amongst other directions equally. The offline dataset contains 150 episodes and
 646 the behavioral policy is ϵ -optimal with noise $\epsilon = 0.1$. Episodes have length $H = 10$.

647 **Sepsis Simulation.** The sepsis simulator [Oberst and Sontag, 2019] is a commonly used envi-
 648 ronment for medically-motivated RL work [Tang and Wiens, 2021]. We use the original authors’
 649 publicly available code: <https://github.com/clinicalml/gumbel-max-scm/tree/sim-v2/sepsisSimDiabetes> (MIT license). The state space consists of five discrete observational vari-
 650 ables (heart rate, blood pressure, oxygen concentration, glucose, diabetes status) and the action
 651 space consists of three different binary treatment options (antibiotic administration, vasopressor ad-
 652 ministration, mechanical ventilation). The probability that each treatment affects the value of each
 653 vital sign is determined by Oberst and Sontag [2019] to reflect patients’ physiology. The ground
 654 truth reward function is sparse and only assigns a positive reward of +1 to surviving patients and
 655 a negative reward of -1 if death occurs (3 or more abnormal vitals) during their stay. The offline
 656 trajectories dataset includes 10,000 episodes following an ϵ -optimal policy with noise $\epsilon = 0.1$ and
 657 the episode length is $H = 20$.
 658

659 D Additional Results

660 We include additional results in this section.

661 In Figure 6, we report the accuracy of the transition and preference model achieved for the Star MDP
 662 as we vary the size of optimality of the offline dataset. Accuracy is measured against all possible
 663 state transitions and over 100 pairs of random trajectories (random combinations of the 5 states and
 664 4 actions in a sequence of $H = 3$). This complements our analysis in Section 7 and fig. 2. We
 665 see a steady improvement in both transition and reward model quality as we increase the amount of
 666 observational data in Figure 6a, which explains the observed dependence of N_p on N_o in Figure 2a.

667 In Figure 6b, we notice low model performance at both extremes of the x-axis. When the dataset is
 668 fully optimal, we find that all trajectories involve the same sequence of actions and states, so learning
 669 a transition or reward model from this data is challenging. We reach a similar conclusion at the other
 670 end of the spectrum at high density ratios, where the coverage the optimal states reduces. We reach

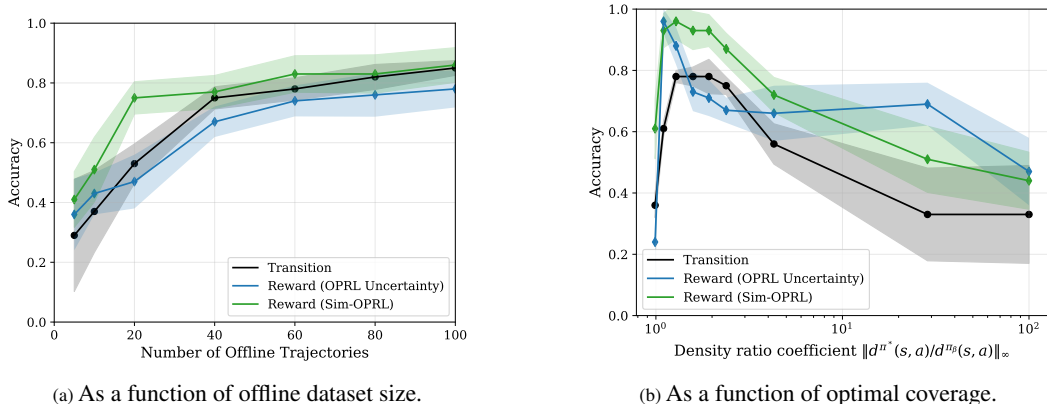


Figure 6: **Transition and preference model accuracy as function of the properties of the observational data** (Star MDP). Preference elicitation is carried out until 10 preferences are queried. Mean and 95% confidence intervals over 20 experiments. Note that the transition model is the same for the two methods, as they have access to the same dataset.

671 highest performance for both models at intermediate values, when diversity of the observational data
 672 is high.

673 Still, it is important to stress that the highest accuracy of both models does not necessarily translate
 674 to the best-performing policy: good performance on the distribution induced by the optimal policy
 675 is more important, as formalized by the concentrability coefficients.

676 Next, we plot performance as a function of preferences sampled for our two additional environments
 677 in Figure 7. We reach similar conclusion to those drawn from the Star MDP in Section 7: within
 678 the offline preference elicitation approaches, OPRL with uniform sampling is the least efficient,
 679 OPRL with uncertainty sampling performs better, and Sim-OPRL even better. The PbOP method
 680 naturally reaches a superior policy with fewer samples as it allows environment interaction and can
 681 thus improve its estimate of the transition model in parallel to learning the preference function.

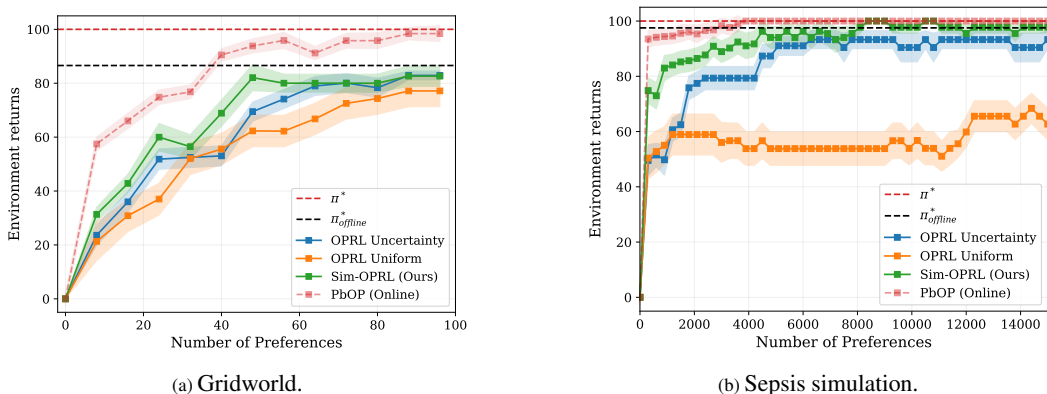


Figure 7: **Empirical results on additional environments.** Mean and 95% confidence interval over 20 experiments. Environment returns are normalised between 0 and 100. Only OPRL and Sim-OPRL are fully offline.

682 E Broader Impact

683 Better preference elicitation strategies for offline reinforcement learning have the potential to facili-
 684 tate and improve decision-making in real-world safety-critical domains like healthcare or economics,
 685 by reducing reliance on direct environment interaction and reducing human effort in providing feed-
 686 back. Potential downsides could include the amplification of biases in the offline data, potentially
 687 leading to suboptimal or unfair policies. Thorough evaluation is therefore crucial to mitigate this

688 before deploying models in such real-world applications. In addition, human preferences may not
689 be fully captured by binary comparisons. As noted in our conclusion, we hope that future work will
690 explore richer feedback mechanisms to better model complex decision-making objectives.