

---

# AMORE: A Model-based Framework for Improving Arbitrary Baseline Policies with Offline Data

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 We propose a new model-based offline RL framework, called Adversarial Models  
2 for Offline Reinforcement Learning (AMORE), which can robustly learn policies  
3 to improve upon an arbitrary baseline policy regardless of data coverage. Based on  
4 the concept of relative pessimism, AMORE is designed to optimize for the worst-  
5 case relative performance when facing uncertainty. In theory, we prove that the  
6 learned policy of AMORE never degrades the performance of the baseline policy  
7 with *any* admissible hyperparameter, and can learn to compete with the best pol-  
8 icy within data coverage when the hyperparameter is well tuned and the baseline  
9 policy is supported by the data. Such a robust policy improvement property makes  
10 AMORE especially suitable for building real-world learning systems, because in  
11 practice ensuring no performance degradation is imperative before considering  
12 any benefit learning can bring.

## 13 1 Introduction

14 Offline reinforcement learning (RL) is a technique for learning decision making policies from logged  
15 data (Jin et al., 2021; Xie et al., 2021a). In comparison with alternate learning techniques, such as  
16 off-policy RL and imitation learning, offline RL reduces the data assumption needed to learn good  
17 policies and does not require collecting new data. Theoretically, offline RL can learn the best policy  
18 that the given data can explain: as long as the offline data includes all scenarios that executing a  
19 near-optimal policy would encounter, an offline RL algorithm can learn a near-optimal policy, even  
20 when the data is collected by highly sub-optimal policies or is not diverse. Such robustness to data  
21 coverage quality makes offline RL a promising technique for solving real-world problems, because  
22 collecting diverse or expert-quality data in practice is expensive or simply infeasible.

23 The fundamental principle behind offline RL is the concept of pessimism in face of uncertainty,  
24 which considers worst-case outcomes for scenarios without data. In implementation, this is realized  
25 by (explicitly or implicitly) constructing performance lower bounds in policy learning, which pe-  
26 nalizes the agent to take uncertain actions. Various designs have been proposed to construct such  
27 lower bounds, including behavior regularization (Fujimoto et al., 2019; Kumar et al., 2019; Wu  
28 et al., 2019; Laroche et al., 2019; Fujimoto and Gu, 2021), point-wise pessimism based on negative  
29 bonuses or truncation (Kidambi et al., 2020; Jin et al., 2021), value penalty (Kumar et al., 2020; Yu  
30 et al., 2020), or two-player games (Cheng et al., 2022; Xie et al., 2021a; Uehara and Sun, 2021).  
31 Conceptually, the tighter the lower bound is, the better the learned policy would perform, as the  
32 performance estimate is more accurate.

33 Despite these advances, offline RL still has not been widely adopted to build learning-based decision  
 34 systems in practice. One reason we posit is that *achieving high performance in the worst case is not*  
 35 *the full picture of designing real-world learning agents.*

36 Usually we apply machine learning to applications that are not completely unknown, but have some  
 37 running policies. These policies are the decision rules that are currently used in the system (e.g.,  
 38 an engineered autonomous driving rule, or a heuristic-based system for diagnosis), and the goal of  
 39 applying a learning algorithm is often to further improve upon these *baseline policies*. As a result,  
 40 it is imperative that the policy learned by the agent does not lead to *performance degradation*. This  
 41 criterion is especially critical for applications where the poor decision outcomes cannot be tolerated  
 42 (such as health care, autonomous driving, and commercial resource allocation).

43 Although optimizing for absolute or relative performance is the same when full information is avail-  
 44 able, they can lead to different policies when we only have partial data coverage. In this case, the  
 45 policy that has the best worst-case performance (which most existing offline RL aims to recover)  
 46 would not necessarily perform better than the baseline policies when deployed in the real envi-  
 47 ronment. Such performance degradation happens when the data do not cover *all* behaviors of the  
 48 baseline policies, which can be due to finite samples or a coverage mismatch between the base-  
 49 lines and the data collection policies. As a result, running policies learned by existing offline RL  
 50 algorithms could risk degrading performance.

51 In this work, we propose a new model-based offline RL framework, called **Adversarial Models for**  
 52 **Offline Reinforcement Learning (AMORE)**, which can robustly learn policies improving upon an  
 53 arbitrary baseline policy. AMORE is designed based on the concept of relative pessimism (Cheng  
 54 et al., 2022), which aims to optimize for the worst-case relative performance when facing uncer-  
 55 tainty. In theory, we prove that the the learned policy from AMORE never degrades the performance  
 56 of the baseline policy of comparison for a wide range of hyperparameters which are given before-  
 57 hand, a property known as Robust Policy Improvement (RPI) (Cheng et al., 2022). In addition, we  
 58 prove that, when the right hyperparameter is chosen and the baseline police is covered by the data,  
 59 the learned policy of AMORE can also compete with any policy within data coverage in an absolute  
 60 sense.

61 To our knowledge, RPI property of offline RL has so far be limited to comparing against the data  
 62 collection policy (Cheng et al., 2022; Fujimoto et al., 2019; Kumar et al., 2019; Wu et al., 2019;  
 63 Laroché et al., 2019; Fujimoto and Gu, 2021). However, it is quite common that the baseline policy  
 64 of interest is different from the data collection policy. For example, in robotics manipulation, we  
 65 often have a dataset of activities different from the target task that we wish to solve. In AMORE, by  
 66 using models, we extend the technique of relative pessimism to achieve RPI with *arbitrary* baseline  
 67 policies, regardless whether the baseline policies collected the data or not.

## 68 2 Preliminaries

69 **Markov Decision Process** We consider an agent acting in an infinite-horizon discounted Markov  
 70 Decision Process (MDP)  $M$  defined by the tuple  $\langle \mathcal{S}, \mathcal{A}, \mathcal{P}, R, \gamma \rangle$  where  $\mathcal{S}$  is the state space,  $\mathcal{A}$   
 71 is the action space,  $\mathcal{P} : \mathcal{S} \times \mathcal{A} \rightarrow \Delta(\mathcal{S})$  is the transition dynamics,  $R : \mathcal{S} \times \mathcal{A} \rightarrow [0, 1]$   
 72 is a scalar reward function and  $\gamma \in [0, 1)$  is the discount factor. The learner selects ac-  
 73 tions using a policy  $\pi : \mathcal{S} \rightarrow \Delta(\mathcal{A})$ . We denote by  $\Pi$  the space of all Markovian poli-  
 74 cies. Let,  $d_M^\pi(s, a)$  denote the discounted state-action distribution obtained by running policy  
 75  $\pi$  on  $M$ , i.e  $d_M^\pi(s, a) = (1 - \gamma) \mathbb{E} [\sum_{t=0}^{\infty} \gamma^t \mathbb{1}(s_t = s, a_t = a | a_t \sim \pi(\cdot | s_t))]$ . Let  $J_M(\pi) =$   
 76  $\mathbb{E}_{\pi, M} [\sum_{t=0}^{\infty} \gamma^t r_t | a_t \sim \pi]$  be the expected discounted return of policy  $\pi$  on  $M$ . The goal of re-  
 77 inforcement learning is to find the policy that maximizes  $J$ . We define the value function as  
 78  $V_M^\pi(s) = \mathbb{E}_{\pi, M} [\sum_{t=0}^{\infty} \gamma^t r_t | s_0 = s]$ , and the related state-action value function (i.e., Q-function)  
 79 as  $Q_M^\pi(s, a) = \mathbb{E}_{\pi, M} [\sum_{t=0}^{\infty} \gamma^t r_t | s_0 = s, s_0 = a]$ . We use  $[0, V_{\max}]$  as the range of value functions.

80 **Offline RL** The aim of offline RL is to output strong policies from a fixed dataset collected using  
 81 a behavior policy without further environmental interactions. We assume the dataset  $\mathcal{D}$  consists of  
 82  $\{(s_i, a_i, r_i, s_{i+1})\}_{i=1}^N$ , where  $(s_i, a_i)$  is sampled i.i.d. from some distribution  $\mu$ . We also abuse  $\mu$

83 as discounted state-action occupancy of behavior policy, i.e.,  $\mu \equiv d_M^\mu$ , and we use  $a \sim \mu(\cdot|s)$  to  
 84 denote sampling from that behavior policy.

85 This paper is concerned with the model-based offline RL problem, and we use  $\mathcal{M}$  to denote the  
 86 model class. For each  $M \in \mathcal{M}$ , we use  $P_M : \mathcal{S} \times \mathcal{A} \rightarrow \Delta(\mathcal{S})$  and  $R_M : \mathcal{S} \times \mathcal{A} \rightarrow [0, 1]$  to denote  
 87 the corresponding transition and reward function of  $M$ .

88 **Assumption 1** (Realizability). *We assume the ground truth model  $M^*$  is in the model class  $\mathcal{M}$ .*

### 89 3 Adversarial Models for Offline Reinforcement Learning (AMORE)

90 In this section, we introduce our proposed approach, Adversarially Trained Models (AMORE), in  
 91 [Algorithm 1](#), and the main theoretical results.

---

#### Algorithm 1 Adversarially Trained Models (AMORE)

---

**Input:** Batch data  $\mathcal{D}$ . Model class  $\mathcal{M}$ . Coefficient  $\alpha$ . Policy class  $\Pi$ . Reference policy  $\pi_{\text{ref}}$ .

1: Construct version space for the model,

$$\mathcal{M}_\alpha = \left\{ M \in \mathcal{M} : \max_{M' \in \mathcal{M}} \mathcal{L}_{\mathcal{D}}(M') - \mathcal{L}_{\mathcal{D}}(M) \leq \alpha \right\}, \quad (1)$$

$$\text{where } \mathcal{L}_{\mathcal{D}}(M) := \sum_{(s,a,r,s') \in \mathcal{D}} \left[ \log \mathbb{P}_M(s'|s, a) - (R_M(s, a) - r)^2 \right], \forall M \in \mathcal{M}. \quad (2)$$

2: Conduct learning via relative pessimism,

$$\hat{\pi} = \operatorname{argmax}_{\pi \in \Pi} \min_{M \in \mathcal{M}_\alpha} J_M(\pi) - J_M(\pi_{\text{ref}}). \quad (3)$$


---

92 AMORE can be viewed as a model-based extension of the ATAC algorithm by [Cheng et al. \(2022\)](#).  
 93 In the next sections, we illustrate that AMORE is not only able to compete with the best data-covered  
 94 policy as prior works (e.g., [Xie et al., 2021a](#); [Uehara and Sun, 2021](#); [Cheng et al., 2022](#)), but also  
 95 enjoys a stronger robust policy improvement guarantee than ([Cheng et al., 2022](#)).

#### 96 3.1 Theoretical Analysis

97 This section analyzes AMORE theoretically and presents guarantees on its absolute performance and  
 98 the policy improvement over the reference policy  $\pi_{\text{ref}}$ . Before presenting the detailed guarantees,  
 99 we introduce generalized single-policy concentrability, which measures the distribution shift over  
 100 some arbitrary policy  $\pi$  and data distribution  $\mu$ .

101 **Definition 1** (Generalized Single-policy Concentrability). *We define the generalized single-policy  
 102 concentrability for policy  $\pi$  for model class  $\mathcal{M}$  and offline data distribution  $\mu$  as*

$$\mathfrak{C}_{\mathcal{M}}(\pi) := \sup_{M \in \mathcal{M}} \frac{\mathbb{E}_{d^\pi} \left[ D_{\text{TV}}(P_M(\cdot|s, a), P^*(\cdot|s, a))^2 + (R_M(s, a) - R^*(s, a))^2 \right]}{d_\mu \left[ D_{\text{TV}}(P_M(\cdot|s, a), P^*(\cdot|s, a))^2 + (R_M(s, a) - R^*(s, a))^2 \right]}.$$

103 Note that  $\mathfrak{C}_{\mathcal{M}}(\pi)$  is always upper bounded by the standard single-policy concentrability coefficient  
 104  $\|d^\pi / \mu\|_\infty$  (e.g., [Jin et al., 2021](#); [Rashidinejad et al., 2021](#); [Xie et al., 2021b](#)), but it can be smaller in  
 105 general with model class  $\mathcal{M}$ . It can also be viewed as a model-based analog of the one in [Xie et al.](#)  
 106 [\(2021a\)](#), and the detailed discussion around  $\mathfrak{C}_{\mathcal{M}}(\pi)$  refers to [Uehara and Sun \(2021\)](#).

107 We are now ready to present the absolute performance guarantee of AMORE.

108 **Theorem 1** (Absolute performance guarantee). *Under [Assumption 1](#), there exists an absolute con-*  
 109 *stant  $c$  such that for any  $\delta \in (0, 1]$ , if we choose  $\alpha = c \cdot (\log(|\mathcal{M}|/\delta))$  in [Algorithm 1](#), then for*  
 110 *arbitrary reference policy  $\pi_{\text{ref}}$  and comparator policy  $\pi^\dagger \in \Pi$ , with probability  $1 - \delta$ , the policy  $\hat{\pi}$*

111 learned by [Algorithm 1](#) satisfies

$$J(\pi^\dagger) - J(\hat{\pi}) \leq \mathcal{O} \left( \left[ \sqrt{\mathfrak{C}_{\mathcal{M}}(\pi^\dagger)} + \sqrt{\mathfrak{C}_{\mathcal{M}}(\pi_{\text{ref}})} \right] \cdot \frac{V_{\max}}{1 - \gamma} \sqrt{\frac{\log(|\mathcal{M}|/\delta)}{n}} \right).$$

112 Roughly speaking, [Theorem 1](#) shows that  $\hat{\pi}$  learned by [Algorithm 1](#) could compete with any policy  
 113  $\pi$  with a large enough dataset, as long as the offline data  $\mu$  has good coverage on comparator policy  
 114  $\pi^\dagger$  (since the reference policy  $\pi_{\text{ref}}$  is the input of [Theorem 1](#), one can set  $\pi_{\text{ref}} = \mu$  (data collection  
 115 policy) as  $\mathfrak{C}_{\mathcal{M}}(\mu) \leq \mathfrak{C}_{\mathcal{M}}(\pi^\dagger)$ ). Compared to the closest model-based offline RL work ([Uehara and  
 116 Sun, 2021](#)), if we set  $\pi_{\text{ref}} = \mu$  (data collection policy), [Theorem 1](#) leads to the almost the same  
 117 guarantee as [Uehara and Sun \(2021, Theorem 1\)](#) (up to constant factors).

118 In addition to the guarantee on the absolute performance above, below we show that, if [Assumption  
 119 1](#) is satisfied and  $\pi_{\text{ref}} \in \Pi$ , AMORE is always guaranteed to improve over  $J(\hat{\pi})$  for a wide range  
 120 choice of pessimistic parameter  $\alpha$ .

121 **Theorem 2** (Robust strong policy improvement). *Under [Assumption 1](#), there exists an absolute  
 122 constant  $c$  such that for any  $\delta \in (0, 1]$ , if: i)  $\alpha \geq c \cdot (\log(|\mathcal{M}|/\delta))$  in [Algorithm 1](#); ii)  $\pi_{\text{ref}} \in \Pi$ , then  
 123 with probability  $1 - \delta$ , the policy  $\hat{\pi}$  learned by [Algorithm 1](#) satisfies  $J(\pi_{\text{ref}}) \geq J(\hat{\pi})$ .*

### 124 3.2 Discussion

125 Improving over some reference policy has been long studied in the literature. To highlight the  
 126 advantage of AMORE, we formally give the definition of different policy improvement properties.

127 **Definition 2** (Robust policy improvement). *Suppose  $\hat{\pi}$  is the learned policy from an algorithm.  
 128 We say the algorithm has the policy improvement (PI) guarantee if  $J(\pi_{\text{ref}}) - J(\hat{\pi}) \leq o(n)/n$  is  
 129 guaranteed for some reference policy  $\pi_{\text{ref}}$  with offline data  $\mathcal{D} \sim \mu$ , where  $n = |\mathcal{D}|$ . We use the  
 130 following two criteria w.r.t.  $\pi_{\text{ref}}$  and  $\mu$  to define different kinds PI:*

- 131 (i) *The PI is strong if  $\pi_{\text{ref}}$  can be selected arbitrarily from policy class  $\Pi$  regardless of the choice  
 132 data-collection policy  $\mu$ ; otherwise, PI is weak (e.g.,  $\pi_{\text{ref}} \equiv \mu$  is required).*
- 133 (ii) *The PI is robust if it can be achieved by a range of hyperparameters with a known subset.*

134 Weak policy improvement is also known as *safe policy improvement* in the literature ([Fujimoto  
 135 et al., 2019](#); [Laroche et al., 2019](#)). It requires the reference policy to be also the behavior policy that  
 136 collects the offline data. In comparison, strong policy improvement imposes a stricter requirement  
 137 on the algorithm, which requires policy improvement *regardless* of how the data were collected.  
 138 This condition is motivated by the common situation where the reference policy is not the data  
 139 collection policy. For example, in a multi-task problem with shared dynamics, the data are collected  
 140 by policies for different tasks, and the reference policy we wish to improve on is task specific. In this  
 141 case, weak policy improvement is meaningless because the behavior policy, which is the average of  
 142 policies from all tasks, does not have meaningful performance in the target task.

143 Since we are learning policies offline, without online interactions, it is not straightforward to tune  
 144 the hyperparameter directly. Therefore, it is desirable that we can design algorithms with these  
 145 properties in a robust manner in terms of hyperparameter selection. Formally, [Definition 2](#) requires  
 146 the policy improvement to be achievable by a set of hyperparameters that is known before learning.

147 [Theorem 2](#) indicates the robust strong policy improvement of AMORE. On the other hand, algo-  
 148 rithms with robust weak policy improvement are available in the literature ([Cheng et al., 2022](#);  
 149 [Fujimoto et al., 2019](#); [Kumar et al., 2019](#); [Wu et al., 2019](#); [Laroche et al., 2019](#); [Fujimoto and Gu,  
 150 2021](#)); this is usually achieved by designing the algorithm to behave like imitation learning (IL)  
 151 for a known set of hyperparameter (e.g., behavior regularization algorithms have a weight that can turn  
 152 off the RL behavior and regress to IL). However, the absolute performance guarantee of achieving  
 153 the best data-covered policy of the IL-like algorithm is challenging due to its imitating nature. To  
 154 our best knowledge, ATAC ([Cheng et al., 2022](#)) is the only algorithm that achieves robust (weak)  
 155 policy improvement as well as guarantees absolute performance.

## 156 References

- 157 Alekh Agarwal, Sham Kakade, Akshay Krishnamurthy, and Wen Sun. Flambe: Structural com-  
158 plexity and representation learning of low rank mdps. *Advances in neural information processing*  
159 *systems*, 33:20095–20107, 2020.
- 160 Ching-An Cheng, Tengyang Xie, Nan Jiang, and Alekh Agarwal. Adversarially trained actor critic  
161 for offline reinforcement learning. *arXiv preprint arXiv:2202.02446*, 2022.
- 162 Scott Fujimoto and Shixiang Shane Gu. A minimalist approach to offline reinforcement learning.  
163 *Advances in neural information processing systems*, 34:20132–20145, 2021.
- 164 Scott Fujimoto, David Meger, and Doina Precup. Off-policy deep reinforcement learning without  
165 exploration. In *International Conference on Machine Learning*, pages 2052–2062, 2019.
- 166 Ying Jin, Zhuoran Yang, and Zhaoran Wang. Is pessimism provably efficient for offline rl? In  
167 *International Conference on Machine Learning*, pages 5084–5096. PMLR, 2021.
- 168 Rahul Kidambi, Aravind Rajeswaran, Praneeth Netrapalli, and Thorsten Joachims. Morel: Model-  
169 based offline reinforcement learning. In *NeurIPS*, 2020.
- 170 Aviral Kumar, Justin Fu, Matthew Soh, George Tucker, and Sergey Levine. Stabilizing off-policy q-  
171 learning via bootstrapping error reduction. *Advances in Neural Information Processing Systems*,  
172 32:11784–11794, 2019.
- 173 Aviral Kumar, Aurick Zhou, George Tucker, and Sergey Levine. Conservative q-learning for offline  
174 reinforcement learning. *Advances in Neural Information Processing Systems*, 33:1179–1191,  
175 2020.
- 176 Romain Laroche, Paul Trichelair, and Remi Tachet Des Combes. Safe policy improvement with  
177 baseline bootstrapping. In *International Conference on Machine Learning*, pages 3652–3661.  
178 PMLR, 2019.
- 179 Qinghua Liu, Alan Chung, Csaba Szepesvári, and Chi Jin. When is partially observable reinforce-  
180 ment learning not scary? In *Conference on Learning Theory*, volume 178, pages 5175–5220.  
181 PMLR, 2022.
- 182 Paria Rashidinejad, Banghua Zhu, Cong Ma, Jiantao Jiao, and Stuart Russell. Bridging offline rein-  
183 forcement learning and imitation learning: A tale of pessimism. *Advances in Neural Information*  
184 *Processing Systems*, 34:11702–11716, 2021.
- 185 Masatoshi Uehara and Wen Sun. Pessimistic model-based offline reinforcement learning under  
186 partial coverage. In *International Conference on Learning Representations*, 2021.
- 187 Sara A van de Geer. *Empirical Processes in M-estimation*, volume 6. Cambridge university press,  
188 2000.
- 189 Yifan Wu, George Tucker, and Ofir Nachum. Behavior regularized offline reinforcement learning.  
190 *arXiv preprint arXiv:1911.11361*, 2019.
- 191 Tengyang Xie, Ching-An Cheng, Nan Jiang, Paul Mineiro, and Alekh Agarwal. Bellman-consistent  
192 pessimism for offline reinforcement learning. *Advances in neural information processing systems*,  
193 34:6683–6694, 2021a.
- 194 Tengyang Xie, Nan Jiang, Huan Wang, Caiming Xiong, and Yu Bai. Policy finetuning: Bridg-  
195 ing sample-efficient offline and online reinforcement learning. *Advances in neural information*  
196 *processing systems*, 34:27395–27407, 2021b.
- 197 Tianhe Yu, Garrett Thomas, Lantao Yu, Stefano Ermon, James Y Zou, Sergey Levine, Chelsea Finn,  
198 and Tengyu Ma. Mopo: Model-based offline policy optimization. *Advances in Neural Information*  
199 *Processing Systems*, 33:14129–14142, 2020.

200 Tong Zhang. From  $\varepsilon$ -entropy to kl-entropy: Analysis of minimum information complexity density  
201 estimation. *The Annals of Statistics*, 34(5):2180–2210, 2006.

# Appendix

202

## 203 A Proofs from Section 3

### 204 A.1 Technical Tools

205 **Lemma 3** (Simulation lemma). *Consider any two MDP model  $M$  and  $M'$ , and any  $\pi : \mathcal{S} \rightarrow \Delta(\mathcal{A})$ ,*  
 206 *we have*

$$|J_M(\pi) - J_{M'}(\pi)| \leq \frac{V_{\max}}{1-\gamma} \mathbb{E}_{d^\pi} [D_{\text{TV}}(P_M(\cdot|s, a), P_{M'}(\cdot|s, a))] + \frac{1}{1-\gamma} \mathbb{E}_{d^\pi} [|R_M(s, a) - R_{M'}(s, a)|].$$

207 **Lemma 3** is the standard simulation lemma in model-based reinforcement learning literature, and its  
 208 proof can be found in, e.g., [Uehara and Sun \(2021, Lemma 7\)](#).

### 209 A.2 Guarantees about Version Space

210 **Lemma 4.** *Let  $M^*$  be the ground truth model. Then, with probability at least  $1 - \delta$ , we have*

$$\max_{M \in \mathcal{M}} \mathcal{L}_{\mathcal{D}}(M) - \mathcal{L}_{\mathcal{D}}(M^*) \leq \mathcal{O}(\log(|\mathcal{M}|/\delta)),$$

211 where  $\mathcal{L}_{\mathcal{D}}$  is defined in [Eq. \(2\)](#).

212 **Proof of Lemma 4.** By [Lemma 6](#), we know

$$\max_{M \in \mathcal{M}} \log \ell_{\mathcal{D}}(M) - \log \ell_{\mathcal{D}}(M^*) \leq \log(|\mathcal{M}|/\delta). \quad (4)$$

213 In addition, by [Xie et al. \(2021a, Theorem A.1\)](#) (with setting  $\gamma = 0$ ), we know w.p.  $1 - \delta$ ,

$$\sum_{(s,a,r,s') \in \mathcal{D}} (R^*(s, a) - r)^2 - \min_{M \in \mathcal{M}} \sum_{(s,a,r,s') \in \mathcal{D}} (R_M(s, a) - r)^2 \lesssim \log(|\mathcal{M}|/\delta). \quad (1)$$

214 Combining the [Eqs. \(1\)](#) and [\(4\)](#), we have w.p.  $1 - \delta$ ,

$$\begin{aligned} & \max_{M \in \mathcal{M}} \mathcal{L}_{\mathcal{D}}(M) - \mathcal{L}_{\mathcal{D}}(M^*) \\ & \leq \max_{M \in \mathcal{M}} \log \ell_{\mathcal{D}}(M) - \min_{M \in \mathcal{M}} \sum_{(s,a,r,s') \in \mathcal{D}} (R_M(s, a) - r)^2 - \mathcal{L}_{\mathcal{D}}(M^*) \\ & \lesssim \log(|\mathcal{M}|/\delta). \end{aligned}$$

215 This completes the proof. □

216 **Lemma 5.** *For any  $M \in \mathcal{M}$ , we have with probability at least  $1 - \delta$ ,*

$$\begin{aligned} & \mathbb{E}_{\mu} \left[ D_{\text{TV}}(P_M(\cdot|s, a), P^*(\cdot|s, a))^2 + (R_M(s, a) - R^*(s, a))^2 \right] \\ & \leq \mathcal{O} \left( \frac{\max_{M' \in \mathcal{M}} \mathcal{L}_{\mathcal{D}}(M') - \mathcal{L}_{\mathcal{D}}(M) + \log(|\mathcal{M}|/\delta)}{n} \right), \end{aligned}$$

217 where  $\mathcal{L}_{\mathcal{D}}$  is defined in [Eq. \(2\)](#).

218 **Proof of Lemma 5.** By [Lemma 7](#), we have w.p.  $1 - \delta$ ,

$$n \cdot \mathbb{E}_{\mu} \left[ D_{\text{TV}}(P_M(\cdot|s, a), P^*(\cdot|s, a))^2 \right] \lesssim \log \ell_{\mathcal{D}}(M^*) - \log \ell_{\mathcal{D}}(M) + \log(|\mathcal{M}|/\delta). \quad (5)$$

219 Also, we have

$$\begin{aligned} & n \cdot \mathbb{E}_{\mu} \left[ (R_M(s, a) - R^*(s, a))^2 \right] \\ & = n \cdot \mathbb{E}_{\mu} \left[ (R_M(s, a) - r)^2 \right] - n \cdot \mathbb{E}_{\mu} \left[ (R^*(s, a) - r)^2 \right] \\ & \quad \text{(see, e.g., [Xie et al., 2021a, Eq. \(A.10\)](#) with  $\gamma = 0$ )} \end{aligned} \quad (6)$$

$$\lesssim \sum_{(s,a,r,s') \in \mathcal{D}} (R_M(s,a) - r)^2 - \sum_{(s,a,r,s') \in \mathcal{D}} (R^*(s,a) - r)^2 + \log(|\mathcal{M}|/\delta),$$

220 where the last inequality is a direct implication of [Xie et al. \(2021a, Lemma A.4\)](#) and  $1 = 1$ .  
 221 Combining [Eqs. \(5\) and \(6\)](#), we obtain

$$\begin{aligned} & n \cdot \mathbb{E}_\mu \left[ D_{\text{TV}}(P_M(\cdot|s,a), P^*(\cdot|s,a))^2 + (R_M(s,a) - R^*(s,a))^2 \right] \\ & \lesssim \log \ell_{\mathcal{D}}(M^*) - \sum_{(s,a,r,s') \in \mathcal{D}} (R^*(s,a) - r)^2 - \log \ell_{\mathcal{D}}(M) + \sum_{(s,a,r,s') \in \mathcal{D}} (R_M(s,a) - r)^2 + \log(|\mathcal{M}|/\delta) \\ & = \mathcal{L}_{\mathcal{D}}(M^*) - \mathcal{L}_{\mathcal{D}}(M) + \log(|\mathcal{M}|/\delta) \\ & \leq \max_{M' \in \mathcal{M}} \mathcal{L}_{\mathcal{D}}(M') - \mathcal{L}_{\mathcal{D}}(M) + \log(|\mathcal{M}|/\delta). \end{aligned}$$

222 This completes the proof.

223

□

### 224 A.3 MLE Guarantees

225 We use  $\ell_{\mathcal{D}}(M)$  to denote the likelihood of model  $M = (P, R)$  with offline data  $\mathcal{D}$ , where

$$\ell_{\mathcal{D}}(M) = \prod_{(s,a,r,s') \in \mathcal{D}} P_M(s'|s,a). \quad (7)$$

226 For the analysis around maximum likelihood estimation, we largely follow the proving idea of [Agarwal et al. \(2020\)](#); [Liu et al. \(2022\)](#), which is inspired by [Zhang \(2006\)](#).  
 227

228 The next lemma shows that the ground truth model  $M^*$  has a comparable log-likelihood compared  
 229 with MLE solution.

230 **Lemma 6.** *Let  $M^*$  be the ground truth model. Then, with probability at least  $1 - \delta$ , we have*

$$\max_{M \in \mathcal{M}} \log \ell_{\mathcal{D}}(M) - \log \ell_{\mathcal{D}}(M^*) \leq \log(|\mathcal{M}|/\delta). \quad (8)$$

231 **Proof of Lemma 6.** The proof of this lemma is obtained by a standard argument of MLE (see, e.g.,  
 232 [van de Geer, 2000](#)). For any  $M \in \mathcal{M}$ ,

$$\begin{aligned} \mathbb{E} [\exp(\log \ell_{\mathcal{D}}(M) - \log \ell_{\mathcal{D}}(M^*))] &= \mathbb{E} \left[ \frac{\ell_{\mathcal{D}}(M)}{\ell_{\mathcal{D}}(M^*)} \right] \\ &= \mathbb{E} \left[ \frac{\prod_{(s,a,r,s') \in \mathcal{D}} \mathbb{P}_M(s'|s,a)}{\prod_{(s,a,r,s') \in \mathcal{D}} \mathbb{P}_{M^*}(s'|s,a)} \right] \\ &= \mathbb{E} \left[ \prod_{(s,a,r,s') \in \mathcal{D}} \frac{\mathbb{P}_M(s'|s,a)}{\mathbb{P}_{M^*}(s'|s,a)} \right] \\ &= \mathbb{E} \left[ \prod_{(s,a) \in \mathcal{D}} \mathbb{E} \left[ \frac{\mathbb{P}_M(s'|s,a)}{\mathbb{P}_{M^*}(s'|s,a)} \mid s, a \right] \right] \\ &= \mathbb{E} \left[ \prod_{(s,a) \in \mathcal{D}} \sum_{s',r} \mathbb{P}_M(s'|s,a) \right] \\ &= 1. \end{aligned} \quad (9)$$

233 Then by Markov's inequality, we obtain

$$\begin{aligned} & \mathbb{P}[(\log \ell_{\mathcal{D}}(M) - \log \ell_{\mathcal{D}}(M^*)) > \log(1/\delta)] \\ & \leq \underbrace{\mathbb{E} [\exp(\log \ell_{\mathcal{D}}(M) - \log \ell_{\mathcal{D}}(M^*))]}_{=1 \text{ by Eq. (9)}} \cdot \exp[-\log(1/\delta)] = \delta. \end{aligned}$$

234 Therefore, taking a union bound over  $\mathcal{M}$ , we obtain

$$\mathbb{P}[(\log \ell_{\mathcal{D}}(M) - \log \ell_{\mathcal{D}}(M^*)) > \log(|\mathcal{M}|/\delta)] \leq \delta.$$

235 This completes the proof.  $\square$

236 The following lemma shows that, the on-support error of any model  $M \in \mathcal{M}$  can be captured via  
237 its log-likelihood (by comparing with the MLE solution).

238 **Lemma 7.** For any  $M = (P, R)$ , we have with probability at least  $1 - \delta$ ,

$$\mathbb{E}_{\mu} \left[ D_{\text{TV}}(P(\cdot|s, a), P^*(\cdot|s, a))^2 \right] \leq \mathcal{O} \left( \frac{\log \ell_{\mathcal{D}}(M^*) - \log \ell_{\mathcal{D}}(M) + \log(|\mathcal{M}|/\delta)}{n} \right),$$

239 where  $\ell_{\mathcal{D}}(\cdot)$  is defined in Eq. (7).

240 **Proof of Lemma 7.** By Agarwal et al. (2020, Lemma 25), we have

$$\mathbb{E}_{\mu} \left[ D_{\text{TV}}(P(\cdot|s, a), P^*(\cdot|s, a))^2 \right] \leq -2 \log \mathbb{E}_{\mu \times P^*} \left[ \exp \left( -\frac{1}{2} \log \left( \frac{P^*(s'|s, a)}{P(s'|s, a)} \right) \right) \right] \quad (10)$$

$$\mathbb{E}_{\mu} \left[ D_{\text{TV}}(R(\cdot|s, a), R^*(\cdot|s, a))^2 \right] \leq -2 \log \mathbb{E}_{\mu \times R^*} \left[ \exp \left( -\frac{1}{2} \log \left( \frac{R^*(r|s, a)}{R(r|s, a)} \right) \right) \right],$$

241 where  $\mu \times P^*$  and  $\mu \times R^*$  denote the ground truth offline joint distribution of  $(s, a, s')$  and  $(s, a, r)$ .

242 Let  $\tilde{\mathcal{D}} = \{(\tilde{s}_i, \tilde{a}_i, \tilde{r}_i, \tilde{s}'_i)\}_{i=1}^n \sim \mu$  be another offline dataset that is independent to  $\mathcal{D}$ . Then,

$$\begin{aligned} & -n \cdot \log \mathbb{E}_{\mu \times P^*} \left[ \exp \left( -\frac{1}{2} \log \left( \frac{P^*(s'|s, a)}{P(s'|s, a)} \right) \right) \right] \\ &= -\sum_{i=1}^n \log \mathbb{E}_{(\tilde{s}_i, \tilde{a}_i, \tilde{s}'_i) \sim \mu} \left[ \exp \left( -\frac{1}{2} \log \left( \frac{P^*(\tilde{s}'_i|\tilde{s}_i, \tilde{a}_i)}{P(\tilde{s}'_i|\tilde{s}_i, \tilde{a}_i)} \right) \right) \right] \\ &= -\log \mathbb{E}_{\tilde{\mathcal{D}} \sim \mu} \left[ \exp \left( \sum_{i=1}^n -\frac{1}{2} \log \left( \frac{P^*(\tilde{s}'_i|\tilde{s}_i, \tilde{a}_i)}{P(\tilde{s}'_i|\tilde{s}_i, \tilde{a}_i)} \right) \right) \middle| \mathcal{D} \right] \\ &= -\log \mathbb{E}_{\tilde{\mathcal{D}} \sim \mu} \left[ \exp \left( \sum_{(s, a, s') \in \tilde{\mathcal{D}}} -\frac{1}{2} \log \left( \frac{P^*(s'|s, a)}{P(s'|s, a)} \right) \right) \middle| \mathcal{D} \right]. \end{aligned} \quad (11)$$

243 We use  $\ell_P(s, a, s')$  as the shorthand of  $-\frac{1}{2} \log \left( \frac{P^*(s'|s, a)}{P(s'|s, a)} \right)$ , for any  $(s, a, s') \in \mathcal{S} \times \mathcal{A} \times \mathcal{S}$ . By Agarwal  
244 et al. (2020, Lemma 24) (see also Liu et al., 2022, Lemma 15), we know

$$\mathbb{E}_{\mathcal{D} \sim \mu} \left[ \exp \left( \sum_{(s, a, s') \in \mathcal{D}} \ell_P(s, a, s') - \log \mathbb{E}_{\tilde{\mathcal{D}} \sim \mu} \left[ \exp \left( \sum_{(s, a, s') \in \tilde{\mathcal{D}}} \ell_P(s, a, s') \right) \middle| \mathcal{D} \right] - \log |\mathcal{M}| \right) \right] \leq 1.$$

245 Thus, we can use Chernoff method as well as a union bound on the equation above to obtain the  
246 following exponential tail bound: with probability at least  $1 - \delta$ , we have for all  $(P, R) = M \in \mathcal{M}$ ,

$$-\log \mathbb{E}_{\tilde{\mathcal{D}} \sim \mu} \left[ \exp \left( \sum_{(s, a, s') \in \tilde{\mathcal{D}}} \ell_P(s, a, s') \right) \middle| \mathcal{D} \right] \leq -\sum_{(s, a, s') \in \mathcal{D}} \ell_P(s, a, s') + 2 \log(|\mathcal{M}|/\delta). \quad (12)$$

247 Plugging back the definition of  $\ell_P$  and combining Eqs. (10) to (12), we obtain

$$n \cdot \mathbb{E}_{\mu} \left[ D_{\text{TV}}(P(\cdot|s, a), P^*(\cdot|s, a))^2 \right] \leq \frac{1}{2} \sum_{(s, a, s') \in \mathcal{D}} \log \left( \frac{P^*(s|s, a)}{P(s'|s, a)} \right) + 2 \log(|\mathcal{M}|/\delta). \quad (13)$$

248 By the same steps of obtaining to Eq. (13), we also have

$$n \cdot \mathbb{E}_{\mu} \left[ D_{\text{TV}}(R(\cdot|s, a), R^*(\cdot|s, a))^2 \right] \leq \frac{1}{2} \sum_{(s, a, r') \in \mathcal{D}} \log \left( \frac{R^*(s|s, a)}{R(s'|s, a)} \right) + 2 \log(|\mathcal{M}|/\delta). \quad (14)$$

249 Combining Eqs. (13) and (14), we obtain

$$\begin{aligned}
& n \cdot \mathbb{E}_\mu \left[ D_{\text{TV}}(P(\cdot|s, a), P^*(\cdot|s, a))^2 + D_{\text{TV}}(R(\cdot|s, a), R^*(\cdot|s, a))^2 \right] \\
& \lesssim \sum_{(s, a, s') \in \mathcal{D}} \log \left( \frac{P^*(s|s, a)}{P(s'|s, a)} \right) + \sum_{(s, a, r') \in \mathcal{D}} \log \left( \frac{R^*(s|s, a)}{R(s'|s, a)} \right) + \log(|\mathcal{M}|/\delta) \\
& = \log \ell_{\mathcal{D}}(M^*) - \log \ell_{\mathcal{D}}(M) + \log(|\mathcal{M}|/\delta). \quad (\ell_{\mathcal{D}}(\cdot) \text{ is defined in Eq. (7)})
\end{aligned}$$

250 This completes the proof.  $\square$

#### 251 A.4 Proof of Main Theorems

252 **Proof of Theorem 1.** By the optimality of  $\hat{\pi}$  (from Eq. (3)), we have

$$\begin{aligned}
J(\pi^\dagger) - J(\hat{\pi}) &= J(\pi^\dagger) - J(\pi_{\text{ref}}) - [J(\hat{\pi}) - J(\pi_{\text{ref}})] \\
&\leq J(\pi^\dagger) - J(\pi_{\text{ref}}) - \min_{M \in \mathcal{M}_\alpha} [J_M(\hat{\pi}) - J_M(\pi_{\text{ref}})] \\
&\quad \text{(by Lemma 6, we have } M^* \in \mathcal{M}_\alpha) \\
&\leq J(\pi^\dagger) - J(\pi_{\text{ref}}) - \min_{M \in \mathcal{M}_\alpha} [J_M(\pi^\dagger) - J_M(\pi_{\text{ref}})], \quad (15)
\end{aligned}$$

253 where the last step is because of  $\pi^\dagger \in \Pi$ . By the simulation lemma (Lemma 3), we know for any  
254 policy  $\pi$  and any  $M \in \mathcal{M}_\alpha$ ,

$$\begin{aligned}
|J(\pi) - J_M(\pi)| &\leq \frac{V_{\max}}{1-\gamma} \mathbb{E}_{d^\pi} [D_{\text{TV}}(P_M(\cdot|s, a), P^*(\cdot|s, a))] + \frac{1}{1-\gamma} \mathbb{E}_{d^\pi} [|R_M(s, a) - R^*(s, a)|] \\
&\leq \frac{V_{\max}}{1-\gamma} \sqrt{\mathbb{E}_{d^\pi} [D_{\text{TV}}(P_M(\cdot|s, a), P^*(\cdot|s, a))^2]} + \frac{1}{1-\gamma} \sqrt{\mathbb{E}_{d^\pi} [(R_M(s, a) - R^*(s, a))^2]} \\
&\lesssim \frac{V_{\max}}{1-\gamma} \sqrt{\mathbb{E}_{d^\pi} [D_{\text{TV}}(P_M(\cdot|s, a), P^*(\cdot|s, a))^2 + (R_M(s, a) - R^*(s, a))^2]} \\
&\quad (a \lesssim b \text{ means } a \leq \mathcal{O}(b)) \\
&\leq \frac{V_{\max} \sqrt{\mathfrak{C}_{\mathcal{M}}(\pi)}}{1-\gamma} \sqrt{\mathbb{E}_\mu [D_{\text{TV}}(P_M(\cdot|s, a), P^*(\cdot|s, a))^2 + (R_M(s, a) - R^*(s, a))^2]} \\
&\lesssim \frac{V_{\max} \sqrt{\mathfrak{C}_{\mathcal{M}}(\pi)}}{1-\gamma} \sqrt{\frac{\max_{M' \in \mathcal{M}} \mathcal{L}_{\mathcal{D}}(M') - \mathcal{L}_{\mathcal{D}}(M) + \log(|\mathcal{M}|/\delta)}{n}} \\
&\quad \text{(by Lemma 5)} \\
&\lesssim \frac{V_{\max} \sqrt{\mathfrak{C}_{\mathcal{M}}(\pi)}}{1-\gamma} \sqrt{\frac{\log(|\mathcal{M}|/\delta)}{n}} \quad (16)
\end{aligned}$$

255 where the last step is because  $\max_{M' \in \mathcal{M}} \mathcal{L}_{\mathcal{D}}(M') - \mathcal{L}_{\mathcal{D}}(M) \leq \alpha = \mathcal{O}(\log(|\mathcal{M}|/\delta)/n)$  by Eq. (1).

256 Combining Eqs. (15) and (16), we obtain

$$J(\pi^\dagger) - J(\hat{\pi}) \lesssim \left[ \sqrt{\mathfrak{C}_{\mathcal{M}}(\pi^\dagger)} + \sqrt{\mathfrak{C}_{\mathcal{M}}(\pi_{\text{ref}})} \right] \cdot \frac{V_{\max}}{1-\gamma} \sqrt{\frac{\log(|\mathcal{M}|/\delta)}{n}}.$$

257 This completes the proof.  $\square$

#### **Proof of Theorem 2.**

$$\begin{aligned}
J(\pi_{\text{ref}}) - J(\hat{\pi}) &= J(\pi_{\text{ref}}) - J(\pi_{\text{ref}}) - [J(\hat{\pi}) - J(\pi_{\text{ref}})] \\
&\leq - \min_{M \in \mathcal{M}_\alpha} [J_M(\hat{\pi}) - J_M(\pi_{\text{ref}})] \quad \text{(by Lemma 6, we have } M^* \in \mathcal{M}_\alpha) \\
&= - \max_{\pi \in \Pi} \min_{M \in \mathcal{M}_\alpha} [J_M(\pi) - J_M(\pi_{\text{ref}})] \quad \text{(by the optimality of } \hat{\pi} \text{ from Eq. (3))} \\
&\leq - \min_{M \in \mathcal{M}_\alpha} [J_M(\pi_{\text{ref}}) - J_M(\pi_{\text{ref}})] \quad (\pi_{\text{ref}} \in \Pi) \\
&= 0.
\end{aligned}$$

258  $\square$