

# LEASE: OFFLINE PREFERENCE-BASED REINFORCEMENT LEARNING WITH HIGH SAMPLE EFFICIENCY

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Offline preference-based reinforcement learning (PbRL) provides an effective way to overcome the challenges of designing reward and the high costs of online interaction. However, since labeling preference needs real-time human feedback, acquiring sufficient preference labels is challenging. To solve this, this paper proposes a **offLine prEference-bAsed RL with high Sample Efficiency (LEASE)** algorithm, where a learned transition model is leveraged to generate unlabeled preference data. Considering the pretrained reward model may generate incorrect labels for unlabeled data, we design an uncertainty-aware mechanism to ensure the performance of reward model, where only high confidence and low variance data are selected. Moreover, we provide the generalization bound of reward model to analyze the factors influencing reward accuracy, and demonstrate that the policy learned by LEASE has theoretical improvement guarantee. The developed theory is based on state-action pair, which can be easily combined with other offline algorithms. The experimental results show that LEASE can achieve comparable performance to baseline under fewer preference data without online interaction.

## 1 INTRODUCTION

Reinforcement learning (RL) has been widely applied in robot control (Haarnoja et al., 2024; Radosavovic et al., 2024) and healthcare (Li et al., 2021; Wu et al., 2022) fields. However, in real-world scenarios, RL faces two serious challenges: 1) It is dangerous for agent to interact directly with environment, especially human-in-loop control (Levine et al., 2020); 2) It is difficult to design a quantitative reward function to accurately reflect the human intention or preference (Hadfield-Menell et al., 2017). Many researches have been explored to solve the above challenges.

Offline RL, learning policy from previous collected dataset without interaction with environment, has become an effective way for challenge 1 (Kumar et al., 2020; Kostrikov et al., 2021a). Preference-based RL (PbRL), which learns the reward function from human preference data without the demand for tedious hand-engineered reward design (Christiano et al., 2017), has been developed recently for challenge 2. However, compared with data  $(s, a, s')$ , the collecting cost of preference data is higher since it demands human real-time feedback to label preference, which often brings high sample complexity (low sample efficiency) (Liang et al., 2022; Park et al., 2022).

Some data augmentation techniques have been proposed to solve problem of the limitation of preference data (Park et al., 2022; Hu et al., 2024). However, the above methods need online interaction with environment and still demand large human feedback. Therefore, this paper aims to design a high sample efficiency offline PbRL algorithm that can achieve comparable control performance with baseline from a limited amount of preference data without interaction with environment.

In offline PbRL setting, existing methods (Shin et al., 2021; Kim et al., 2023; Zhang et al., 2024) typically involve two steps: reward learning and policy learning. However, the theoretical guarantees of reward generalization and policy improvement are not provided. Therefore, this paper focuses on answering the following two questions: 1) In algorithm, how to use the limited preference data to learn accurate reward model and guarantee agent performance? 2) In theory, what factors affect the reward generalization and what is the detailed relationship between offline RL algorithm itself, the reward gap and the improvement of policy performance?

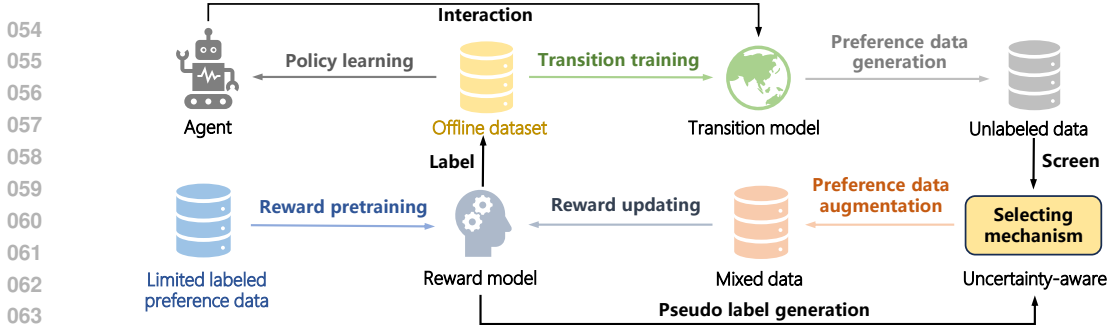


Figure 1: The learning framework of LEASE. The offline dataset without reward is used to train transition model and limited labeled preference dataset is used to pretrain reward model. Then the generated unlabeled preference data is screened through uncertainty-aware selecting mechanism. The reward model is updated based on labeled and generated dataset. Finally, the agent learns the policy based on offline dataset and learned reward model.

For question 1, motivated by model-based offline RL, which learns the environment model to broaden data coverage (Yu et al., 2020; Liu et al., 2024), we also train the transition model to achieve data augmentation (improve sample efficiency). Specifically, the agent generates two different trajectories through interaction with the learned transition model, and then the pretrained reward model generates pseudo label for these two trajectories. However, the preference model may generate incorrect pseudo label, which may bring unstable training and low performance.

For question 2, there are very fewer algorithms for offline PbRL theory. The most relevant work is Zhan et al. (2024), where the theory of policy improvement guarantee is established. However, the generalization bound for reward model is not considered, and this theory is based on the whole trajectory rather than state-action pair in offline algorithms, which is not conducive to the theoretical analysis of specific offline PbRL algorithms since many offline RL are based on state-action pair.

To solve the above problems, this paper proposes a novel offLine prEference-bAsed RL with high Sample Efficiency algorithm (LEASE), where a selecting mechanism is designed to guarantee the quantity of generated dataset. Moreover, we develop a new generalization bound for reward model, and provide the theory of policy improvement guarantee for offline PbRL based on state-action pair. Fig. 1 shows the learning framework of LEASE. LEASE can achieve comparable performance to baseline under fewer preference data. The contributions are given below:

- 1) A novel learning framework, named LEASE, is proposed for offline PbRL, where an innovative preference augmentation technique is utilized to improve sample efficiency.
- 2) An uncertainty-aware mechanism is designed for screening the generated preference data so as to guarantee the stability of reward training and improve the accuracy of reward prediction.
- 3) The generalization bound of reward model and the theory of policy improvement based on state-action pair are developed to analyze factors that influence reward and policy performance.

In this paper, the theoretical and experimental results show that LEASE has policy improvement guarantee and can achieve superior performance under fewer preference data on D4RL benchmark.

## 2 PRELIMINARIES

**Offline Reinforcement Learning.** The framework of RL is based on the Markov Decision Process (MDP) that is described by the tuple  $\mathcal{M} = (\mathcal{S}, \mathcal{A}, R, T, \rho, \gamma)$ , where  $\mathcal{S}$  is the state space,  $\mathcal{A}$  is the action space,  $R : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$  is the reward function,  $T : \mathcal{S} \times \mathcal{A} \rightarrow \Delta(\mathcal{S})$  is the transition dynamics,  $\rho$  is the initial state distribution, and  $\gamma \in (0, 1)$  is the discount factor (Levine et al., 2020). The term  $\Delta(\Omega)$  denotes the set of probability distribution over space  $\Omega$ . The goal of RL is to optimize the policy  $\pi$  that maximizes the expected discounted return  $J(\pi, R) := \mathbb{E}_{(s,a) \sim d_T^\pi(s,a)} [R(s,a)] / (1-\gamma)$ , where  $d_T^\pi(s,a) := d_T^\pi(s)\pi(a|s)$  is the state-action marginal distribution under the learned policy  $\pi$ . The discounted state marginal distribution  $d_T^\pi(s)$  is denoted as  $(1-\gamma) \sum_{t=0}^{\infty} \gamma^t \mathcal{P}(s_t = s | \pi)$ ,

where  $\mathcal{P}(s_t = s|\pi)$  is the probability of reaching state  $s$  at time  $t$  by rolling out  $\pi$ . The policy  $\pi$  can be derived from  $Q$ -learning, which learns the  $Q$ -function that satisfies Bellman operator:  $\mathcal{T}Q(s, a) := R(s, a) + \gamma \mathbb{E}_{s' \sim \mathcal{T}_{\mathcal{M}}(s'|s, a)}[\max_{a' \in \mathcal{A}} Q(s', a')]$  (Sutton et al., 1998).

The goal of offline RL is to learn a policy from offline dataset  $\mathcal{D}_{\text{offline}} = \{(s_i, a_i, r_i, s'_i)\}_{i=1}^C$  collected by the behavior policy  $\mu$ . The policy learning includes two parts: policy evaluation (minimizing the Bellman error) and policy improvement (maximizing the  $Q$ -function), that is

$$\begin{aligned} \hat{Q} &\leftarrow \arg \min_Q \mathbb{E}_{s, a, s' \sim \mathcal{D}} [Q(s, a) - \hat{\mathcal{T}}\hat{Q}(s, a)]^2 && \text{(Policy Evaluation)} \\ \hat{\pi} &\leftarrow \arg \max_{\pi} \mathbb{E}_{s \sim \mathcal{D}, a \sim \pi} [\hat{Q}(s, a)] && \text{(Policy Improvement)} \end{aligned} \quad (1)$$

where  $\mathcal{D}$  can be a fixed offline dataset or replay buffer generated by the current policy  $\hat{\pi}$  interacting with the environment. The operator  $\hat{\mathcal{T}}$  is the Bellman operator based on sample, denoted as  $\hat{\mathcal{T}}Q(s, a) := R(s, a) + \gamma \max_{a' \in \mathcal{A}} Q(s', a')$  (Kumar et al., 2020).

**Preference-based Reinforcement Learning.** Different from standard RL setting, the reward function is not available in PbRL. Instead, PbRL learns the reward function  $\hat{R}$  from preferences between pairs of trajectory segments to align the human intention (Wilson et al., 2012), where a trajectory segment  $\sigma$  of length  $L$  is defined as a sequence of states and actions  $\{s_k, a_k, \dots, s_{k+L-1}, a_{k+L-1}\} \in (\mathcal{S} \times \mathcal{A})^L$ . Given a pair of segments  $(\sigma_0, \sigma_1)$ , human choose which segment is preferred, i.e.,  $y \in \{0, 1\}$  (Christiano et al., 2017). The preference label  $y = 1$  indicates  $\sigma_1 \succ \sigma_0$  and  $y = 0$  indicates  $\sigma_0 \succ \sigma_1$ , where  $\sigma_i \succ \sigma_j$  denotes that the segment  $i$  is preferable to the segment  $j$ .

The preference data is stored in the dataset  $\mathcal{D}_l$  as a tuple  $(\sigma_0, \sigma_1, y)$ . To obtain the reward function  $\hat{R}$  parameterized by  $\psi$ , the Bradley-Terry model (Bradley and Terry, 1952) is utilized to define a preference predictor following previous works (Kim et al., 2023; Verma and Metcalfe, 2024):

$$P(\sigma_0 \succ \sigma_1; \psi) = \frac{\exp \sum_t \hat{R}_{\psi}(s_t^0, a_t^0)}{\exp \sum_t \hat{R}_{\psi}(s_t^0, a_t^0) + \exp \sum_t \hat{R}_{\psi}(s_t^1, a_t^1)}. \quad (2)$$

Then, based on the collected preference dataset  $\mathcal{D}_l$ , the reward function  $\hat{R}_{\psi}$  can be updated through minimizing the cross-entropy loss between the predicted and the true preference labels (Verma and Metcalfe, 2024):

$$\mathcal{L}_R(\psi) = -\mathbb{E}_{(\sigma_0, \sigma_1, y) \sim \mathcal{D}_l} [(1 - y) \log P(\sigma_0 \succ \sigma_1; \psi) + y \log P(\sigma_1 \succ \sigma_0; \psi)]. \quad (3)$$

After obtaining the reward function  $\hat{R}$ , offline PbRL optimizes the policy by maximizing the expected discounted return  $J(\pi, R)$  like standard offline RL.

**Model-based Offline Reinforcement Learning.** Model-based RL learns the dynamics model to improve sample efficiency (Rigter et al., 2022; Liu et al., 2024). They used the offline dataset  $\mathcal{D}_{\text{offline}}$  to estimate transition model  $\hat{T}(s'|s, a)$ . The transition model  $\hat{T}_{\varphi}$  parameterized by  $\varphi$  is typically trained via maximum likelihood estimation (MLE):

$$\mathcal{L}_T(\varphi) = -\mathbb{E}_{(s, a, s') \sim \mathcal{D}_{\text{offline}}} [\log \hat{T}_{\varphi}(s' | s, a)]. \quad (4)$$

In practical implantation, the transition model is approximated by Gaussian distribution and the MLE method is employed to train  $N_T$  ensemble transition models  $\{T_{\varphi}^i = \mathcal{N}(\mu_{\varphi}^i, \sigma_{\varphi}^i)\}_{i=1}^{N_T}$  (Yu et al., 2020). The samples are generated through  $H$ -step rollouts. Here, we use the trained transition to generate more trajectory to achieve data augmentation.

### 3 PROBLEM FORMULATION

The cost of collecting preference data is high since it needs real-time human feedback. Therefore, this paper aims to improve sample efficiency to learn accurate reward model from limited preference data and guarantee agent performance. For PbRL, the ideal form is that the learned reward function  $\hat{R}$  from collected preference data can be consistent with true reward  $R^*$ . Here, we define the function class as  $\mathcal{R} = \{R : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}\}$ . Then we assume the reward class  $\mathcal{R}$  is realizable.

**Assumption 1 (Realizability).** Let  $d(s, a) \in \Delta(\mathcal{S} \times \mathcal{A})$  be the arbitrary data distribution. Then, for any distribution  $d(s, a)$ ,  $\inf_{R \in \mathcal{R}} \mathbb{E}_{(s, a) \sim d(s, a)} [R^*(s, a) - R(s, a)]^2 < \varepsilon_R$  holds.

Based on the above assumption, the optimal reward function can be obtained if the preference data can cover all data space. However, preference data tend to be limited, even less than hundred in real-world scenarios, such as rehabilitation robot field. Therefore, different from previous PbRL (Park et al., 2022; Hu et al., 2024; Zhang et al., 2024), this paper learns reward function from the limited dataset  $\mathcal{D}_l = \{(\sigma_0^l, \sigma_1^l, y)^{(i)}\}_{i=1}^{N_l}$ . We train the transition model through Eq. (4) to generate more unlabeled preference data  $\mathcal{D}_u = \{(\sigma_0^u, \sigma_1^u)^{(i)}\}_{i=1}^{N_u}$ . The pseudo labels  $\hat{y}$  for unlabeled dataset  $\mathcal{D}_u$  can be obtained through reward model trained on  $\mathcal{D}_l$  (Park et al., 2022), that is

$$\hat{y}(\sigma_0^u, \sigma_1^u) = \mathbb{1}[P(\sigma_1^u \succ \sigma_0^u; \psi) > 0.5], \quad (5)$$

where  $\mathbb{1}(\cdot)$  is indicator function. The reward function is updated through collected labeled dataset  $\mathcal{D}_l$  and generated unlabeled dataset  $\mathcal{D}_u$ . Then, the reward model can be optimized through minimizing

$$\mathcal{L}'_R(\psi) = \mathbb{E}_{(\sigma_0^l, \sigma_1^l, y) \sim \mathcal{D}_l} [L((\sigma_0^l, \sigma_1^l), y)] + \mathbb{E}_{(\sigma_0^u, \sigma_1^u) \sim \mathcal{D}_u} [L((\sigma_0^u, \sigma_1^u), \hat{y})], \quad (6)$$

where  $L((\sigma_0, \sigma_1), y) = -(1-y) \log P(\sigma_0 \succ \sigma_1; \psi) - y \log P(\sigma_1 \succ \sigma_0; \psi)$ . However, the pretrained reward model may generate incorrect pseudo-labels for unlabeled dataset, leading to noisy training and poor generalization (Rizve et al., 2021). Therefore, one key question is how to design a data selecting mechanism to improve prediction accuracy and guarantee training stability.

The another key aspect is the theory for offline PbRL. There are very few algorithms specifically designed for offline PbRL with strong guarantee, including generalization bound for reward model and safe improvement guarantee for policy learning. Therefore, another key question is to develop a systematic theory for offline PbRL, including generalization bound and improvement guarantee.

## 4 REWARD LEARNING

The reward learning involves two stages: pretraining based on collected labeled data  $\mathcal{D}_l$  and updating based on  $\mathcal{D}_l$  and unlabeled data  $\mathcal{D}_u$  during policy learning, which is essentially semi-supervised learning (Berthelot et al., 2019). This section focuses on designing a data selecting function  $f(\sigma_0^u, \sigma_1^u)$  to ensure the quality of generated preference data, and explaining the factors that influence the generalization ability of reward model.

### 4.1 UNCERTAINTY-AWARE PSEUDO-LABEL SELECTION FOR REWARD LEARNING

Motivated by the previous pseudo-labeling work (Rizve et al., 2021), we select data from unlabeled dataset  $\mathcal{D}_u$  according two principles: confidence and uncertainty. The data with high confidence and low uncertainty can be chosen for reward training. High confidence refers that pre-trained reward model discriminates the preference of two trajectories into  $\hat{y}$  with high probability  $p(\sigma_0^u, \sigma_1^u, \hat{y})$ . Low uncertainty refers that the  $N_R$  reward models (model ensembles) predicts with small variance  $\tau(\sigma_0^u, \sigma_1^u, N_R)$ . The probability confidence  $p(\sigma_0^u, \sigma_1^u, \hat{y})$  and uncertainty variance  $\tau(\sigma_0^u, \sigma_1^u, N_R)$  can be denoted as

$$\begin{aligned} p(\sigma_0^u, \sigma_1^u, \hat{y}) &= (1 - \hat{y}) \cdot \bar{P}(\sigma_0^u \succ \sigma_1^u; \psi) + \hat{y} \cdot \bar{P}(\sigma_1^u \succ \sigma_0^u; \psi) \\ \tau(\sigma_0^u, \sigma_1^u, N_R) &= \text{Std} \{ \mathbf{P}(\sigma_0^u \succ \sigma_1^u; \psi) \}, \end{aligned} \quad (7)$$

where  $N_R$  is the number of reward model,  $\text{Std} \{ \mathbf{P}(\cdot) \}$  denotes the variance of output probability of  $N_R$  reward models, and  $\bar{P}(\cdot)$  is the mean probability of  $N_R$  reward models. Therefore, according to high confidence and low uncertainty, the  $f(\sigma_0^u, \sigma_1^u)$  can be denoted as:

$$f(\sigma_0^u, \sigma_1^u) = \mathbb{1} \left[ p(\sigma_0^u, \sigma_1^u, \hat{y}) > \kappa_p \right] \cdot \mathbb{1} \left[ \tau(\sigma_0^u, \sigma_1^u, N_R) < \kappa_\tau \right], \quad (8)$$

where  $\kappa_p$  and  $\kappa_\tau$  are thresholds of confidence and uncertainty respectively.  $f(\sigma_0^u, \sigma_1^u) = 1$  denotes the generated data  $(\sigma_0^u, \sigma_1^u)$  is selected, and  $f(\sigma_0^u, \sigma_1^u) = 0$  denotes the data is not selected. Then combing Eq. (6), the reward model can be optimized through minimizing

$$\hat{\mathcal{L}}_R(\psi) = \frac{1}{N_l} \sum_{i=1}^{N_l} L((\sigma_0^l, \sigma_1^l)^{(i)}, y^i) + \frac{1}{\tilde{N}_u} \sum_{j=1}^{N_u} f((\sigma_0^u, \sigma_1^u)^{(j)}) \cdot L((\sigma_0^u, \sigma_1^u)^{(j)}, \hat{y}^j), \quad (9)$$

where  $\tilde{N}_u$  is the number of generated dataset after screening, and  $\hat{\mathcal{L}}_R(\psi)$  is the empirical risk. It contains two parts: labeled loss  $\hat{\mathcal{L}}_l(\psi)$  and unlabeled loss  $\hat{\mathcal{L}}_u(\psi)$ . Note that the label in unlabeled

loss is pseudo-label, which may be different from the true label. Therefore, there is the gap between the unlabeled loss with pseudo-label  $\widehat{\mathcal{L}}_u(\psi)$  and that with true label  $\widehat{\mathcal{L}}'_u(\psi)$ . Before bounding this gap, we firstly give the below assumption without loss of generality.

**Assumption 2.** For the pretrained reward model  $\widehat{R}_\psi$  through the limited labeled dataset  $\mathcal{D}_l = \{(\sigma_0^l, \sigma_1^l, y)^{(i)}\}_{i=1}^{N_l}$ , if the pseudo label  $\widehat{y}$  is defined in Eq. (5), then the pseudo-labeling error for the unlabeled dataset  $\mathcal{D}_u = \{(\sigma_0^u, \sigma_1^u)^{(j)}\}_{j=1}^{N_u}$  is smaller than  $\eta$ , i.e.,  $\sum_{j=1}^{N_u} \mathbb{1}[\widehat{y}^j \neq y^j]/N_u \leq \eta$ .

**Proposition 1.** Suppose that the loss  $L((\sigma_0, \sigma_1), y)$  is bounded by  $\Omega$ . Then, for any  $R \in \mathcal{R}$ , under Assumption 2 the following equation holds:

$$|\widehat{\mathcal{L}}_u(\psi) - \widehat{\mathcal{L}}'_u(\psi)| \leq \eta\Omega. \quad (10)$$

The proof of Proposition 1 can be found in Appendix A.1. This gap is mainly influenced by the term  $\eta$ . Through selecting high confidence and low uncertainty samples, the  $\eta$  can be significantly reduced. Combining Eqs. (5) and (8), we find that if  $p(\sigma_0^u, \sigma_1^u, \widehat{y})$  is small, the  $L((\sigma_0^u, \sigma_1^u), \widehat{y})$  would become larger. Too large error can lead to unstable training and low performance. Therefore, through selecting mechanism  $f(\sigma_0^u, \sigma_1^u)$ , a more accurate subset of pseudo-labels can be used in reward training, and through setting a higher confidence threshold for pseudo labeling, a lower prediction error and the stability of training can be guaranteed.

#### 4.2 GENERALIZATION BOUND FOR REWARD MODEL

Before developing generalization bound, we define the expected error  $\mathcal{L}_R(\psi)$  with respect to the reward model  $R(s, a; \psi)$ :  $\mathcal{L}_R(\psi) = \mathbb{E}_{(\sigma_0, \sigma_1, y) \sim \mathcal{D}}[L((\sigma_0, \sigma_1), y)]$ , where  $\mathcal{D}$  is data distribution. The reward model is trained through minimizing the empirical error  $\widehat{\mathcal{L}}_R(\psi)$  in Eq. (9). Through developing generalization error bound, we can analyse the factors that influence generalization ability.

Similar to the previous generalization bound works (Xie et al., 2024), Rademacher complexity, which measures the richness of a certain hypothesis space (Mohri and Muñoz Medina, 2012), is introduced firstly. The definition is given below:

**Definition 1** (Empirical Rademacher complexity). Let  $\mathcal{F}$  be a family of functions mapping from  $\mathcal{Z}$  to  $\mathbb{R}$  and  $\widehat{\mathcal{S}} = \{z_1, \dots, z_{N_s}\}$  be a fixed sample of size  $N_s$  drawn from the distribution  $\mathcal{S}$  over  $\mathcal{Z}$ . The empirical Rademacher complexity of  $\mathcal{G}$  for sample  $\widehat{\mathcal{S}}$  is defined as

$$\widehat{\mathfrak{R}}_{\widehat{\mathcal{S}}}(\mathcal{F}) = \mathbb{E}_{\sigma} \left[ \sup_{f \in \mathcal{F}} \frac{1}{N_s} \sum_{i=1}^{N_s} \sigma_i f(z_i) \right], \quad (11)$$

where  $\sigma = (\sigma_1, \dots, \sigma_{N_s})$  are independent uniform random variables taking values in  $\{-1, +1\}$ .

**Definition 2** (Induced Reward Function Families). Given for a space  $\mathcal{R}$  of reward function  $R$ , the induced reward function families  $\Pi(\mathcal{F})$  is defined as

$$\Pi(\mathcal{R}) = \{(\sigma_0, \sigma_1) \rightarrow -\log P(\sigma_0 \succ \sigma_1) \mid R \in \mathcal{R}\}, \quad (12)$$

where  $P(\cdot)$  is defined in Eq. (2) and  $\Pi(\mathcal{R})$  is the union of projections of  $\mathcal{R}$  onto each dimension.

In general, lower Rademacher complexity corresponds to better generalization performance. Then, based on the above definition and Proposition 1, we develop the new generalization bound for reward model trained through Eq. (9). The specific theorem is given below.

**Theorem 1.** Let reward model be trained on the labeled dataset  $\mathcal{D}_l = \{(\sigma_0^l, \sigma_1^l, y)^{(i)}\}_{i=1}^{N_l}$  and unlabeled dataset  $\mathcal{D}_u = \{(\sigma_0^u, \sigma_1^u)^{(i)}\}_{i=1}^{N_u}$ . Then, for any  $\delta > 0$ , with probability at least  $1 - \delta$ , under Assumption 1 and 2 the following holds for any reward function  $R \in \mathcal{R}$ ,

$$\mathcal{L}_R(\psi) \leq \widehat{\mathcal{L}}_R(\psi) + \eta\Omega + 4\widehat{\mathfrak{R}}_{\widehat{\mathcal{D}}}(\Pi(\mathcal{R})) + 3\sqrt{\frac{\log(2/\delta)}{2(N_l + N_u)}}, \quad (13)$$

where  $\widehat{\mathcal{D}}$  is the input combination of labeled and unlabeled dataset, i.e.  $\{(\sigma_0, \sigma_1)^{(i)}\}_{i=1}^{N_l + N_u}$ .

**Algorithm 1** LEASE: offLine prEference-bAsed RL with high Sample Efficiency

**Input:** offline dataset  $\mathcal{D}_{\text{offline}}$ , limited labeled preference dataset  $\mathcal{D}_l$ , critic  $Q_\omega$ , policy  $\pi_\theta$ , transition model  $T_\varphi$ , and reward model  $R_\psi$ .

**Output:** the reward model  $\hat{R}$  and the policy  $\hat{\pi}$  learned by LEASE.

- 1: **Initialization:** Randomly initializing all networks and generated unlabeled dataset  $\mathcal{D}_u = \emptyset$ .
- 2: Train an ensemble transition model  $\{T_\varphi^i\}_{i=1}^{N_T}$  on  $\mathcal{D}_{\text{offline}}$  according to Eq. (4).
- 3: Pre-train an ensemble reward model  $\{R_\psi^i\}_{i=1}^{N_R}$  on  $\mathcal{D}_l$  according to Eq. (3).
- 4: **for**  $t = 1, 2, \dots, N_{\text{iter}}$  **do**
- 5:   **if**  $t \% \text{rollout frequency} = 0$  **then**
- 6:     Generate synthetic  $H$ -step rollouts by  $T_\varphi$ . Add transition trajectory into  $\mathcal{D}_u$ .
- 7:   **end if**
- 8:   **if** *reward update condition* is TRUE **then**
- 9:     Give pseudo label for unlabeled dataset  $\mathcal{D}_u$  based on pre-trained reward model.
- 10:    Select high confidence and low uncertainty unlabeled data according to Eq. (8).
- 11:    Train the reward model based on  $\mathcal{D}_l$  and  $\mathcal{D}_u$  according to Eq. (9).
- 12:   **end if**
- 13:   Train policy using offline RL algorithms, such as CQL and IQL.
- 14: **end for**

The proof of Theorem 1 can be found in Appendix A.2. Theorem 1 indicates that the expected error  $\mathcal{L}_R(\psi)$  is bounded by empirical error  $\hat{\mathcal{L}}_R(\psi)$ , pseudo-labeling error  $\eta\Omega$ , Rademacher complexity and constant terms. The constant term can become small when the number of unlabeled data increases, but it may cause unstable training when labels are inaccurate. According to Proposition 1, empirical error and pseudo-labeling error can be reduced through selecting mechanism  $f(\sigma_0^u, \sigma_1^u)$ , thus the better generalization ability, that is tighter upper bound, can be achieved. Note that Theorem 1 has generality and can be applicable to methods that train reward model using pseudo-labels. The difference among various methods may lie in how they train reward model to improve the accuracy of pseudo-labels, that is reducing  $\eta$ . If one method fails to reduce the pseudo-labeling error  $\eta$ , the upper generalization bound would be looser than that of no data augmentation.

## 5 POLICY LEARNING

This paper directly uses the proposed offline RL algorithms, such as CQL (Kumar et al., 2020) and IQL (Kostrikov et al., 2021b), to perform policy learning. To solve the problem of high cost for labeling preference, we propose a offLine prEference-bAsed RL with high Sample Efficiency algorithm (LEASE). This section aims to describe the details of LEASE and establishes the theory for safe policy improvement guarantee.

### 5.1 THE IMPLEMENTATION DETAILS OF LEASE

Most offline RL algorithms are based on Eq. (1). They proposed various algorithms to solve the problem of distribution shift. In offline PbRL setting, the reward is unknown for agent. The accuracy of reward model directly influence the performance of policy. Moreover, the unlabeled data are generated through the current policy interaction with learned transition model. Therefore, the reward model and policy influence each other. In practical implementation, the policy,  $Q$ -function, reward and transition model are parameterized by  $\pi_\theta$ ,  $Q_\omega$ ,  $R_\psi$  and  $T_\varphi$  respectively.

Algorithm 1 gives pseudocode for LEASE. The goal of LEASE is to learn better reward model  $R_\psi$  from fewer labeled preference dataset  $\mathcal{D}_l$  and learn better policy  $\pi_\theta$  from offline dataset  $\mathcal{D}_{\text{offline}}$  without interaction with environment. The key of LEASE is how to train reward model in policy learning process. Notably, the reward model only update once instead of updating constantly in this process. The reward update condition is set to when the number of unlabeled data in buffer reaches the maximum buffer capacity. The update time is influenced by rollout frequency and rollout batch size. The reward update condition, rollout length  $H$  and selecting mechanism  $f(\sigma_1, \sigma_2)$  are the important factors for reward model. The more detailed analysis can be found in Appendix B.3.

## 5.2 SAFE POLICY IMPROVEMENT GUARANTEE

Offline RL aims to guarantee  $\xi$ -safe policy improvement over the behavior policy  $\mu$  (the policy used to collect offline dataset), *i.e.*  $J(\hat{\pi}, R^*) \geq J(\mu, R^*) - \xi$  (Kumar et al., 2020). This part develops theoretical guarantee of policy improvement for LEASE, that is giving the bound for  $J(\mu, R^*) - J(\hat{\pi}, R^*)$ , and further analyzes sample complexity when reward model is unknown. Firstly, we propose a new single-policy concentrability coefficient for PbRL based on state-action pair instead of trajectory like (Zhan et al., 2024).

**Definition 3** (Concentrability coefficient for PbRL). *Concentrability coefficient  $\mathcal{C}_{\mathcal{R}}(\pi)$  is used to measure how well reward model errors transfer between offline data distribution  $d_T^\mu$  and the visitation distribution  $d_T^\pi$  under transition  $T$  and policy  $\pi$ , defined as*

$$\mathcal{C}_{\mathcal{R}}(\pi) = \sup_{\hat{R} \in \mathcal{R}} \frac{\left| \mathbb{E}_{(s,a) \sim d_T^\pi} [R^*(s,a) - \hat{R}(s,a)] \right|}{\left| \mathbb{E}_{(s,a) \sim d_T^\mu} [R^*(s,a) - \hat{R}(s,a)] \right|}. \quad (14)$$

where  $R^*(s, a)$  is the true reward model and coefficient  $\mathcal{C}_{\mathcal{R}}(\pi)$  is upper bounded by  $\|d_T^\pi/d_T^\mu\|_\infty$ .

Next, similar to (Geer, 2000; Zhan et al., 2024), we use  $\varepsilon$ -bracketing number to measure the complexity of reward function class  $\mathcal{R}$ , which can be defined as

**Definition 4** ( $\varepsilon$ -bracketing number). *The  $\varepsilon$ -bracketing number  $\mathcal{N}_{\mathcal{R}}(\varepsilon)$  is the minimum number of  $\varepsilon$ -brackets  $(l, u)$  required to cover a function class  $\mathcal{R}$ , where each bracket  $(l, u)$  satisfies  $l(\sigma_0, \sigma_1) \leq P_R(\sigma_0, \sigma_1) \leq u(\sigma_0, \sigma_1)$  and  $\|l - u\|_1 \leq \varepsilon$  for all  $R \in \mathcal{R}$  and all trajectory-pairs  $(\sigma_0, \sigma_1)$ , and  $P_R(\sigma_0, \sigma_1)$  is the probability that segment  $\sigma_0$  is preferable to segment  $\sigma_1$  defined in Eq. (2).*

**Theorem 2.** *Under Assumption 1, for any  $\delta \in (0, 1]$ , the policy  $\hat{\pi}$  learned by LEASE, with high probability  $1 - \delta$ , satisfies that  $J(\mu, R^*) - J(\hat{\pi}, R^*)$  is upper bound by*

$$\xi + \frac{1 + \mathcal{C}_{\mathcal{R}}(\hat{\pi})}{1 - \gamma} \left( \sqrt{\frac{4C}{NL^2} \log \left( \frac{\mathcal{N}_{\mathcal{R}}(1/N)}{\delta} \right)} + \sqrt{\frac{4R_{max}^2 \log(1/\delta)}{NL}} \right), \quad (15)$$

where  $C > 0$  is the absolute constant defined in Eq. (A.19),  $N$  is the size of preference dataset, and  $L$  is the length of trajectory. The term  $\xi$  is performance gap depending on offline algorithm itself.

The proof of Theorem 2 can be found in Appendix A.3. The  $\xi$  is constant when offline algorithm is determined. Therefore, the tighter bound can be achieved through reducing performance gap caused by reward model (see Proposition 3). It can be reduced to small value  $\varpi$  with sample complexity of

$$N = \tilde{O} \left( \frac{4(1 + \mathcal{C}_{\mathcal{R}}(\hat{\pi}))^2}{\varpi^2(1 - \gamma)^2} \left( \sqrt{\frac{C \log(\mathcal{N}_{\mathcal{R}}(1/N)/\delta)}{L^2}} + \sqrt{\frac{R_{max}^2 \log(1/\delta)}{L}} \right)^2 \right), \quad (16)$$

where concentrability coefficient  $\mathcal{C}_{\mathcal{R}}(\hat{\pi})$  can become smaller when the learned policy  $\hat{\pi}$  is close to behavior policy  $\mu$ . Bracketing number  $\log(\mathcal{N}_{\mathcal{R}}(1/N))$  measures the complexity of function class  $\mathcal{R}$  and takes  $\tilde{O}(d)$  in linear reward model (Zhan et al., 2024). Notably, LEASE can learn accurate reward model through data augmentation under fewer  $N_l$  preference data ( $N_l < N$ ). Thus, the performance gap can become tighter under fewer data, that is LEASE can have higher sample efficiency compared with Zhan et al. (2024). More discussion can be found in Appendix C.1.

## 6 EXPERIMENTS

This section focuses on answering the following questions: **Q<sub>1</sub>**: How well the accuracy and generalization of reward model through data augmentation under fewer preference data? (Section 6.2) **Q<sub>2</sub>**: How well does LEASE perform compared with offline PbRL baseline in standard benchmark tasks? (Section 6.1) We answer these questions using D4RL benchmark (Fu et al., 2020) with several control tasks. More hyperparameters and implementation details are provided in Appendix B. The code for LEASE is available at [github.com/\\*\\*/](https://github.com/**/).

Table 1: Results for the D4RL tasks during the last 5 iterations of training averaged over 3 seeds.  $\pm$  captures the standard deviation over seeds. Bold indicates the performance within 2% of the best performing algorithm for each offline algorithm.

| Task Name             | CQL (Kumar et al., 2020)          |                                  |                                   | IQL (Kostrikov et al., 2021b)     |                                  |                                   |
|-----------------------|-----------------------------------|----------------------------------|-----------------------------------|-----------------------------------|----------------------------------|-----------------------------------|
|                       | URLHF                             | FEWER                            | LEASE                             | URLHF                             | FEWER                            | LEASE                             |
| walker2d-m            | 76.0 $\pm$ 0.9                    | 77.4 $\pm$ 0.6                   | <b>78.4 <math>\pm</math> 0.9</b>  | <b>78.4 <math>\pm</math> 0.5</b>  | 71.8 $\pm$ 0.8                   | 74.6 $\pm$ 1.8                    |
| walker2d-m-e          | 92.8 $\pm$ 22.4                   | 77.7 $\pm$ 0.3                   | <b>98.6 <math>\pm</math> 18.1</b> | <b>109.4 <math>\pm</math> 0.1</b> | 105.2 $\pm$ 3.1                  | <b>108.1 <math>\pm</math> 0.5</b> |
| hopper-m              | 54.7 $\pm$ 3.4                    | 55.8 $\pm$ 2.8                   | <b>56.5 <math>\pm</math> 0.6</b>  | 50.8 $\pm$ 4.2                    | <b>56.7 <math>\pm</math> 2.6</b> | <b>56.0 <math>\pm</math> 0.5</b>  |
| hopper-m-e            | <b>57.4 <math>\pm</math> 4.9</b>  | 53.6 $\pm$ 0.9                   | 56.4 $\pm$ 0.8                    | <b>94.3 <math>\pm</math> 7.4</b>  | 56.0 $\pm$ 1.4                   | 55.9 $\pm$ 1.9                    |
| halfcheetah-m         | <b>43.4 <math>\pm</math> 0.1</b>  | <b>43.5 <math>\pm</math> 0.1</b> | <b>43.5 <math>\pm</math> 0.1</b>  | <b>43.3 <math>\pm</math> 0.2</b>  | 42.2 $\pm$ 0.1                   | <b>43.0 <math>\pm</math> 0.3</b>  |
| halfcheetah-m-e       | <b>62.7 <math>\pm</math> 7.1</b>  | 48.3 $\pm$ 0.7                   | 53.2 $\pm$ 3.1                    | <b>91.0 <math>\pm</math> 2.3</b>  | 59.1 $\pm$ 4.9                   | 62.4 $\pm$ 1.4                    |
| <b>Mujoco Average</b> | <b>64.5 <math>\pm</math> 6.5</b>  | 59.4 $\pm$ 0.9                   | <b>64.4 <math>\pm</math> 4.0</b>  | <b>77.9 <math>\pm</math> 2.5</b>  | 65.2 $\pm$ 2.2                   | 66.7 $\pm$ 1.0                    |
| pen-human             | <b>9.8 <math>\pm</math> 14.1</b>  | 0.5 $\pm$ 3.0                    | 3.8 $\pm$ 4.6                     | 50.2 $\pm$ 15.8                   | 67.3 $\pm$ 10.0                  | <b>75.6 <math>\pm</math> 3.3</b>  |
| pen-expert            | <b>138.3 <math>\pm</math> 5.2</b> | 128.1 $\pm$ 0.7                  | 132.5 $\pm$ 2.3                   | <b>132.9 <math>\pm</math> 4.6</b> | 104.1 $\pm$ 12.9                 | 113.8 $\pm$ 6.3                   |
| door-human            | <b>4.7 <math>\pm</math> 5.9</b>   | 0.2 $\pm$ 1.0                    | <b>4.7 <math>\pm</math> 8.8</b>   | 3.5 $\pm$ 3.2                     | 4.0 $\pm$ 2.5                    | <b>5.9 <math>\pm</math> 0.5</b>   |
| door-expert           | <b>103.9 <math>\pm</math> 0.8</b> | 103.0 $\pm$ 0.9                  | <b>103.2 <math>\pm</math> 0.7</b> | <b>105.4 <math>\pm</math> 0.4</b> | 104.5 $\pm$ 0.6                  | <b>105.2 <math>\pm</math> 0.2</b> |
| hammer-human          | <b>0.9 <math>\pm</math> 0.3</b>   | 0.3 $\pm$ 0.0                    | 0.3 $\pm$ 0.0                     | 1.4 $\pm$ 1.0                     | 1.2 $\pm$ 0.2                    | <b>1.7 <math>\pm</math> 0.4</b>   |
| hammer-expert         | 120.2 $\pm$ 6.8                   | 124.1 $\pm$ 2.1                  | <b>126.3 <math>\pm</math> 1.2</b> | <b>127.4 <math>\pm</math> 0.2</b> | 125.2 $\pm$ 2.3                  | <b>126.3 <math>\pm</math> 0.1</b> |
| <b>Adroit Average</b> | <b>63.0 <math>\pm</math> 5.5</b>  | 59.4 $\pm$ 1.3                   | 61.8 $\pm$ 3.0                    | <b>70.1 <math>\pm</math> 4.2</b>  | 67.1 $\pm$ 4.8                   | <b>71.4 <math>\pm</math> 1.8</b>  |

## 6.1 RESULTS ON BENCHMARK TASKS

This part chooses several mujoco and adroit tasks to evaluate the performance of LEASE, where offline dataset and preference dataset originate from Fu et al. (2020) and Yuan et al. (2024), respectively. Yuan et al. (2024) provided baseline for offline PbRL based on 2000 preference data labeled by two types (We denote this method as URLHF). One is crowd-sourced labels obtained by crowd-sourcing, and the other is synthetic labels based on ground truth reward. To test LEASE performance for two types preference data, the labels of preference is from human-realistic feedback for mujoco tasks, and ground truth reward for adroit tasks. This paper aims to test the performance under fewer preference data without online interaction. We take the first 100 data of (Yuan et al., 2024) as 2000 preference dataset, and denote method without data augmentation as FEWER. The effects of the number of labeled preference data  $N_l$  can be found in Appendix B.3.

Table 2: The comparison results between the performance using selecting mechanism and that not using. The latter method is denoted as FRESH.  $\uparrow$  denotes the improvement of performance.

| Task Name  | walker2d-m                          | hopper-m               | halfcheetah-m          | walker2d-m-e            | hopper-m-e             | halfcheetah-m-e         |
|------------|-------------------------------------|------------------------|------------------------|-------------------------|------------------------|-------------------------|
| <b>CQL</b> | <b>FRESH</b> 76.0 $\pm$ 1.3         | 54.4 $\pm$ 1.6         | 43.4 $\pm$ 0.1         | 77.4 $\pm$ 2.2          | 55.0 $\pm$ 0.4         | 49.9 $\pm$ 0.9          |
|            | <b>LEASE</b> 78.4 $\uparrow$ (3.3%) | 56.5 $\uparrow$ (4.0%) | 43.5 $\uparrow$ (0.6%) | 98.6 $\uparrow$ (27.4%) | 56.4 $\uparrow$ (2.6%) | 53.2 $\uparrow$ (6.6%)  |
| <b>IQL</b> | <b>FRESH</b> 72.7 $\pm$ 4.1         | 53.5 $\pm$ 0.8         | 42.5 $\pm$ 0.1         | 105.3 $\pm$ 1.0         | 53.9 $\pm$ 0.1         | 54.5 $\pm$ 2.2          |
|            | <b>LEASE</b> 74.6 $\uparrow$ (2.6%) | 56.0 $\uparrow$ (4.5%) | 43.0 $\uparrow$ (1.3%) | 108.1 $\uparrow$ (2.7%) | 55.9 $\uparrow$ (3.7%) | 62.4 $\uparrow$ (14.6%) |

Here, we evaluate LEASE performance based on CQL (Kumar et al., 2020) and IQL (Kostrikov et al., 2021b) two offline algorithms. Table 1 gives the results of the average normalized score with standard deviation, where m and m-e denote medium and medium-expert respectively. The offline PbRL algorithm under sufficient feedback (URLHF) are obtained from paper (Yuan et al., 2024). This table shows that LEASE can greatly improve performance under fewer preference data compared with FEWER, and achieve comparable performance to URLHF that demand large human feedback. This validates the content of Theorem 2 that LEASE can reduce the performance gap caused by reward model under fewer preference data and has high sample efficiency. Table 2 shows the effect of designed selecting mechanism for agent performance, where the method not using selecting mechanism is denoted as FRESH. This table indicates that the application of selecting mechanism  $f(\sigma_0, \sigma_1)$  can efficiently improve agent performance. The comparison results between LEASE and URLHF under fewer preference data, and results of model-based offline RL algorithm under the designed framework can be found in Appendix B.4.

## 6.2 RESULTS FOR REWARD MODEL

Learning an accurate reward is difficult in the offline PbRL setting when preference labels are scarce. Fig. 2 illustrates the relationship between the prediction accuracy of preference model without data augmentation and the number of preference data  $N_l$ , where preference model is based on reward



432  
433  
434  
435  
436  
437  
438  
439  
440  
441  
442  
443



Figure 3: The comparison between prediction value by the learned rewards and their ground truths for different methods under (a) hopper-medium-expert and (b) halfcheetah-medium-expert datasets, where the value predicted by the trained reward model and ground-truth reward value are both normalized to  $[0, 1]$ . From left to right are methods LEASE, FEWER and FRESH, respectively.

444  
445  
446  
447  
448  
449  
450  
451  
452  
453  
454  
455

Table 3: The performance of the learned transition model, where error is calculated by the sum of the mean square values of predicted and true value in each dimension.

| Task Name | walker2d-m       | walker2d-m-e    | hopper-m        | hopper-m-e      | halfcheetah-m    | halfcheetah-m-e  |
|-----------|------------------|-----------------|-----------------|-----------------|------------------|------------------|
| Error     | $10.19 \pm 0.30$ | $6.57 \pm 0.10$ | $0.35 \pm 0.01$ | $0.38 \pm 0.01$ | $16.18 \pm 0.21$ | $16.19 \pm 0.48$ |
| Task Name | pen-human        | pen-expert      | door-human      | door-expert     | hammer-human     | hammer-expert    |
| Error     | $1.02 \pm 0.08$  | $1.05 \pm 0.01$ | $0.06 \pm 0.04$ | $0.04 \pm 0.01$ | $0.21 \pm 0.05$  | $0.58 \pm 0.03$  |

456  
457  
458  
459  
460  
461  
462  
463  
464  
465  
466  
467  
468  
469  
470  
471  
472  
473  
474  
475  
476

model (Eq. (2)). Here, we use the last 200 of 2000 preference data that are not seen in training stage as evaluation dataset. It indicates that as the number of preference data  $N_l$  increases, the prediction accuracy of the preference model shows an upward trend. This also validates Theorem 1 that increasing the number of samples  $N_l$  can reduce the generalization error.

Next, we compare the accuracy of reward model before and after data augmentation under fewer preference dataset. Fig. 3 shows the comparison results between prediction and ground truth of reward, where the predicted and true rewards are both normalized to  $[0, 1]$ . We randomly sample 500 data from unlabeled datasets that are not seen in training stage for evaluation. This figure indicates the linear relationship between reward predicted by LEASE and ground truth is better than that of other two methods. The prediction of reward model learned by method FRESH is very narrow and the accuracy is greatly reduced compared with FEWER. This is because generated data of FRESH are not selected by  $f(\sigma_0, \sigma_1)$ , which causes substantial errors for the labels of generated data and leads to the collapse of training. This validates that the performance of reward model can be improved through data augmentation and selecting mechanism, where the designed selecting mechanism has a greater impact on reward model performance. The more related experimental results can be found in Appendix B.4.

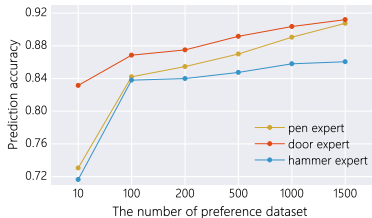


Figure 2: The relationship between preference model accuracy and the number of preference dataset  $N_l$ .

477  
478  
479  
480  
481  
482  
483  
484  
485

### 6.3 RESULTS FOR TRANSITION MODEL

The accuracy of transition model is the key factor influencing agent performance. Table 3 shows the accuracy result for the learned transition model, where the accuracy of transition model is measured by the sum of the mean square values of the predicted value and the true value in each dimension. Note that, for walker2d and halfcheetah tasks, the error of the learned transition model is greatly larger than hopper tasks in mujoco environment. This is mainly because the state space of walker2d and hopper is larger than hopper, and the physical model of them are more complex. Intuitively, the lower the accuracy of the transition model, the poorer performance of the agent. The results of the accuracy of the learned transition model for agent performance can be found in Appendix B.4.

## 7 RELATED WORK

**The algorithm for offline PbRL.** Offline PbRL eliminates the demand for interaction with environment and handcraft designed reward. OPAL (Shin et al., 2021) is the first algorithm that combining offline RL and PbRL. PT (Kim et al., 2023) utilized transformer-based architecture to design preference model capable of generating non-Markovian rewards. OPPO (Kang et al., 2023) directly optimized the policy in a high-level embedding space without learning a separate reward function. However, the above methods demands a large amount of preference dataset. Zhang et al. (2024) trains a diffusion model to achieve data augmentation, but this consumers larger training cost. Our method LEASE can achieve superior performance under fewer dataset and time.

**The theory for offline PbRL.** There are few algorithms that provide theoretical guarantees for offline PbRL, including the generalization bound of reward model and the guarantee of policy improvement. Zhu et al. (2023) studied offline PbRL, but the analysis are restricted to linear model. Zhan et al. (2024) extended to general function approximation, but the generalization analysis of reward model is not provided and the theory is based on trajectory. LEASE gives the theoretical analysis for reward model, and provides the theoretical guarantee for policy based on state-action pair. The theory of LEASE can be easily combined with other offline RL theory.

**Semi-supervised learning.** The goal of semi-supervised learning is using unlabeled data to improve model performance when labeled data is limited. Data selection techniques is used to filter the data with clean labels from a noisy dataset. Han et al. (2018) selected unlabeled data with small losses from one network. Li et al. (2020) used a two-component GMM to separate the dataset into a clean set and a noisy set. Rizve et al. (2021) modeled the prediction uncertainty of unlabeled data to screen data. Xiao et al. (2023) selected data with high confidence as clean data. Motivated by (Rizve et al., 2021) and (Xiao et al., 2023), we trained ensemble reward model to select data with high confidence and low variance to guarantee the quality of unlabeled data. The more related works for this study can be found in Appendix C.2.

## 8 CONCLUSION

This paper proposes a novel offline PbRL algorithm (LEASE) with high sample efficiency. LEASE can achieve comparable performance under fewer preference dataset. By selecting high confidence and low variance data, the stability and accuracy of reward model are guaranteed. Moreover, this paper provides the theoretical analysis for LEASE, including generalization of reward model and policy improvement guarantee. This theory can be easily connected with other offline algorithms. The theoretical and experimental results demonstrate that the data selecting mechanism  $f(\sigma_0, \sigma_1)$  can effectively improve performance of reward model and the performance learned by LEASE can be guaranteed under fewer preference dataset. However, there is still the gap between the true and the learned reward model. Future works can focus on how to further reduce this gap under the limited preference data or how to achieve conservative estimation for state-action pairs where the learned reward model predicts inaccurately.

## REFERENCES

- Gaon An, Seungyong Moon, Jang-Hyun Kim, and Hyun Oh Song. Uncertainty-based offline reinforcement learning with diversified q-ensemble. *Advances in Neural Information Processing Systems*, 34:7436–7447, 2021.
- David Berthelot, Nicholas Carlini, Ian Goodfellow, Nicolas Papernot, Avital Oliver, and Colin A Raffel. Mixmatch: A holistic approach to semi-supervised learning. *Advances in neural information processing systems*, 32, 2019.
- Ralph Allan Bradley and Milton E Terry. Rank analysis of incomplete block designs: I. the method of paired comparisons. *Biometrika*, 39(3/4):324–345, 1952.
- Paul F Christiano, Jan Leike, Tom Brown, Miljan Martic, Shane Legg, and Dario Amodei. Deep reinforcement learning from human preferences. *Advances in neural information processing systems*, 30, 2017.

- 540 Justin Fu, Aviral Kumar, Ofir Nachum, George Tucker, and Sergey Levine. D4rl: Datasets for deep  
541 data-driven reinforcement learning. *arXiv preprint arXiv:2004.07219*, 2020.
- 542
- 543 Sara A Geer. *Empirical Processes in M-estimation*, volume 6. Cambridge university press, 2000.
- 544
- 545 Tuomas Haarnoja, Ben Moran, Guy Lever, Sandy H Huang, Dhruva Tirumala, Jan Humplik, Markus  
546 Wulfmeier, Saran Tunyasuvunakool, Noah Y Siegel, Roland Hafner, et al. Learning agile soccer  
547 skills for a bipedal robot with deep reinforcement learning. *Science Robotics*, 9(89):eadi8022,  
548 2024.
- 549 Dylan Hadfield-Menell, Smitha Milli, Pieter Abbeel, Stuart J Russell, and Anca Dragan. Inverse  
550 reward design. *Advances in neural information processing systems*, 30, 2017.
- 551
- 552 Bo Han, Quanming Yao, Xingrui Yu, Gang Niu, Miao Xu, Weihua Hu, Ivor Tsang, and Masashi  
553 Sugiyama. Co-teaching: Robust training of deep neural networks with extremely noisy labels.  
554 *Advances in neural information processing systems*, 31, 2018.
- 555
- 556 Wassily Hoeffding. Probability inequalities for sums of bounded random variables. *The collected  
557 works of Wassily Hoeffding*, pages 409–426, 1994.
- 558
- 559 Xiao Hu, Jianxiong Li, Xianyuan Zhan, Qing-Shan Jia, and Ya-Qin Zhang. Query-policy misalign-  
560 ment in preference-based reinforcement learning. In *12th International Conference on Learning  
561 Representations*, 2024.
- 562
- 563 Yachen Kang, Diyuan Shi, Jinxin Liu, Li He, and Donglin Wang. Beyond reward: offline preference-  
564 guided policy optimization. In *Proceedings of the 40th International Conference on Machine  
565 Learning*, pages 15753–15768, 2023.
- 566
- 567 Changyeon Kim, Jongjin Park, Jinwoo Shin, Honglak Lee, Pieter Abbeel, and Kimin Lee. Prefer-  
568 ence transformer: Modeling human preferences using transformers for rl. In *11th International  
569 Conference on Learning Representations*, 2023.
- 570
- 571 Ilya Kostrikov, Rob Fergus, Jonathan Tompson, and Ofir Nachum. Offline reinforcement learning  
572 with fisher divergence critic regularization. In *International Conference on Machine Learning*,  
573 pages 5774–5783. PMLR, 2021a.
- 574
- 575 Ilya Kostrikov, Ashvin Nair, and Sergey Levine. Offline reinforcement learning with implicit q-  
576 learning. In *9th International Conference on Learning Representations*, 2021b.
- 577
- 578 Aviral Kumar, Justin Fu, Matthew Soh, George Tucker, and Sergey Levine. Stabilizing off-policy q-  
579 learning via bootstrapping error reduction. *Advances in Neural Information Processing Systems*,  
580 32, 2019.
- 581
- 582 Aviral Kumar, Aurick Zhou, George Tucker, and Sergey Levine. Conservative q-learning for offline  
583 reinforcement learning. *Advances in Neural Information Processing Systems*, 33:1179–1191,  
584 2020.
- 585
- 586 Sergey Levine, Aviral Kumar, George Tucker, and Justin Fu. Offline reinforcement learning: Tuto-  
587 rial, review, and perspectives on open problems. *arXiv preprint arXiv:2005.01643*, 2020.
- 588
- 589 Junnan Li, Richard Socher, and Steven CH Hoi. Dividemix: Learning with noisy labels as semi-  
590 supervised learning. In *8th International Conference on Learning Representations*, 2020.
- 591
- 592 Minhan Li, Yue Wen, Xiang Gao, Jennie Si, and He Huang. Toward expedited impedance tuning  
593 of a robotic prosthesis for personalized gait assistance by reinforcement learning control. *IEEE  
Transactions on Robotics*, 38(1):407–420, 2021.
- 589
- 590 Xinran Liang, Katherine Shu, Kimin Lee, and Pieter Abbeel. Reward uncertainty for exploration in  
591 preference-based reinforcement learning. In *10th International Conference on Learning Repre-  
592 sentations*, 2022.
- 593
- Qinghua Liu, Alan Chung, Csaba Szepesvári, and Chi Jin. When is partially observable reinforce-  
ment learning not scary? In *Conference on Learning Theory*, pages 5175–5220. PMLR, 2022.

- 594 Xiao-Yin Liu, Xiao-Hu Zhou, Xiao-Liang Xie, Shi-Qi Liu, Zhen-Qiu Feng, Hao Li, Mei-Jiang Gui,  
595 Tian-Yu Xiang, De-Xing Huang, and Zeng-Guang Hou. Domain: Mildly conservative model-  
596 based offline reinforcement learning. *arXiv preprint arXiv:2309.08925*, 2023.  
597
- 598 Xiao-Yin Liu, Xiao-Hu Zhou, Guo-Tao Li, Hao Li, Mei-Jiang Gui, Tian-Yu Xiang, De-Xing Huang,  
599 and Zeng-Guang Hou. Micro: Model-based offline reinforcement learning with a conservative  
600 bellman operator. In *33th International Joint Conference on Artificial Intelligence*, 2024.
- 601 Mehryar Mohri and Andres Muñoz Medina. New analysis and algorithm for learning with drift-  
602 ing distributions. In *International Conference on Algorithmic Learning Theory*, pages 124–138.  
603 Springer, 2012.  
604
- 605 Mehryar Mohri, Afshin Rostamizadeh, and Amreet Talwalkar. *Foundations of machine learning*.  
606 MIT press, 2018.
- 607 Jongjin Park, Younggyo Seo, Jinwoo Shin, Honglak Lee, Pieter Abbeel, and Kimin Lee. Surf:  
608 Semi-supervised reward learning with data augmentation for feedback-efficient preference-based  
609 reinforcement learning. In *10th International Conference on Learning Representations, ICLR*  
610 *2022*. International Conference on Learning Representations, 2022.  
611
- 612 Ilija Radosavovic, Tete Xiao, Bike Zhang, Trevor Darrell, Jitendra Malik, and Koushil Sreenath.  
613 Real-world humanoid locomotion with reinforcement learning. *Science Robotics*, 9(89):eadi9579,  
614 2024.
- 615 Marc Rigter, Bruno Lacerda, and Nick Hawes. Rambo-rl: Robust adversarial model-based offline  
616 reinforcement learning. *Advances in neural information processing systems*, 35:16082–16097,  
617 2022.  
618
- 619 Mamshad Nayeem Rizve, Kevin Duarte, Yogesh S Rawat, and Mubarak Shah. In defense of pseudo-  
620 labeling: An uncertainty-aware pseudo-label selection framework for semi-supervised learning.  
621 In *9th International Conference on Learning Representations*, 2021.
- 622 Daniel Shin, Daniel S Brown, and Anca D Dragan. Offline preference-based apprenticeship learning.  
623 *arXiv preprint arXiv:2107.09251*, 2021.  
624
- 625 Yihao Sun, Jiaji Zhang, Chengxing Jia, Haoxin Lin, Junyin Ye, and Yang Yu. Model-bellman  
626 inconsistency for model-based offline reinforcement learning. In *International Conference on*  
627 *Machine Learning*, pages 33177–33194. PMLR, 2023.
- 628
- 629 Richard S Sutton, Andrew G Barto, et al. Introduction to reinforcement learning. vol. 135, 1998.
- 630 Mudit Verma and Katherine Metcalf. Hindsight priors for reward learning from human preferences.  
631 In *12th International Conference on Learning Representations*, 2024.  
632
- 633 Aaron Wilson, Alan Fern, and Prasad Tadepalli. A bayesian approach for policy learning from  
634 trajectory preference queries. *Advances in neural information processing systems*, 25, 2012.
- 635 Ruofan Wu, Junmin Zhong, Brent Wallace, Xiang Gao, He Huang, and Jennie Si. Human-robotic  
636 prosthesis as collaborating agents for symmetrical walking. *Advances in Neural Information*  
637 *Processing Systems*, 35:27306–27320, 2022.  
638
- 639 Ruixuan Xiao, Yiwen Dong, Haobo Wang, Lei Feng, Runze Wu, Gang Chen, and Junbo Zhao.  
640 Promix: combating label noise via maximizing clean sample utility. In *Proceedings of the Thirty-*  
641 *Second International Joint Conference on Artificial Intelligence*, pages 4442–4450, 2023.
- 642 Ming-Kun Xie, Jiahao Xiao, Hao-Zhe Liu, Gang Niu, Masashi Sugiyama, and Sheng-Jun Huang.  
643 Class-distribution-aware pseudo-labeling for semi-supervised multi-label learning. *Advances in*  
644 *Neural Information Processing Systems*, 36, 2024.  
645
- 646 Tianhe Yu, Garrett Thomas, Lantao Yu, Stefano Ermon, James Y Zou, Sergey Levine, Chelsea Finn,  
647 and Tengyu Ma. Mopo: Model-based offline policy optimization. *Advances in Neural Information*  
*Processing Systems*, 33:14129–14142, 2020.

648 Tianhe Yu, Aviral Kumar, Rafael Rafailov, Aravind Rajeswaran, Sergey Levine, and Chelsea Finn.  
649 Combo: Conservative offline model-based policy optimization. *Advances in Neural Information*  
650 *Processing Systems*, 34:28954–28967, 2021.  
651

652 Yifu Yuan, HAO Jianye, Yi Ma, Zibin Dong, Hebin Liang, Jinyi Liu, Zhixin Feng, Kai Zhao, and  
653 YAN ZHENG. Uni-rlhf: Universal platform and benchmark suite for reinforcement learning with  
654 diverse human feedback. In *12th International Conference on Learning Representations*, 2024.

655 Wenhao Zhan, Masatoshi Uehara, Nathan Kallus, Jason D Lee, and Wen Sun. Provable offline  
656 preference-based reinforcement learning. In *12th International Conference on Learning Repre-*  
657 *sentations*, 2024.

658 Zhilong Zhang, Yihao Sun, Junyin Ye, Tian-Shuo Liu, Jiaji Zhang, and Yang Yu. Flow to bet-  
659 ter: Offline preference-based reinforcement learning via preferred trajectory generation. In *12th*  
660 *International Conference on Learning Representations*, 2024.  
661

662 Banghua Zhu, Michael Jordan, and Jiantao Jiao. Principled reinforcement learning with human feed-  
663 back from pairwise or k-wise comparisons. In *International Conference on Machine Learning*,  
664 pages 43037–43067. PMLR, 2023.  
665  
666  
667  
668  
669  
670  
671  
672  
673  
674  
675  
676  
677  
678  
679  
680  
681  
682  
683  
684  
685  
686  
687  
688  
689  
690  
691  
692  
693  
694  
695  
696  
697  
698  
699  
700  
701

## A RELATED PROOFS

This section gives the detailed proof for the theorems in the main text.

### A.1 PROOF OF PROPOSITION 1

**The Proof for Proposition 1:** This proof follows the previous work (Xie et al., 2024). Obviously, the largest pseudo-labeling error  $\eta$  in Assumption 2 includes exactly two types of pseudo-labeling error:

$$\begin{aligned}\eta_1 &= \frac{1}{N_u} \sum_{j=1}^{N_u} \mathbb{1} \left[ P \left( \sigma_1^j \succ \sigma_0^j; \psi \right) > 0.5, y^j = 0 \right], \\ \eta_2 &= \frac{1}{N_u} \sum_{j=1}^{N_u} \mathbb{1} \left[ P \left( \sigma_0^j \succ \sigma_1^j; \psi \right) > 0.5, y^j = 1 \right],\end{aligned}\tag{A.1}$$

where  $\eta_1$  represents the error ratio of classifying the first category as the second category, and  $\eta_2$  is the error ratio of classifying the second category as the first category. Then, we prove the gap between the unlabeled loss with pseudo-label  $\widehat{\mathcal{L}}_u(\psi)$  and that with true label  $\widehat{\mathcal{L}}'_u(\psi)$  from two sides:  $\widehat{\mathcal{L}}_u(\psi) \leq \widehat{\mathcal{L}}'_u(\psi) + \eta\Omega$  and  $\widehat{\mathcal{L}}_u(\psi) \geq \widehat{\mathcal{L}}'_u(\psi) - \eta\Omega$ . Notably, we ignore the selecting mechanism in the below proof since it doesn't influence the proof result.

**Step 1:** Proving upper bound :  $\widehat{\mathcal{L}}_u(\psi) \leq \widehat{\mathcal{L}}'_u(\psi) + \eta\Omega$ .

$$\begin{aligned}\widehat{\mathcal{L}}_u(\psi) &= \frac{1}{N_u} \sum_{j=1}^{N_u} L \left( (\sigma_0^u, \sigma_1^u)^{(j)}, \widehat{y}^j \right) \\ &= -\frac{1}{N_u} \sum_{j=1}^{N_u} \mathbb{1} \left[ P \left( \sigma_0^j \succ \sigma_1^j; \psi \right) > 0.5 \right] \log P \left( \sigma_0^j \succ \sigma_1^j; \psi \right) \\ &\quad + \mathbb{1} \left[ P \left( \sigma_1^j \succ \sigma_0^j; \psi \right) > 0.5 \right] \log P \left( \sigma_1^j \succ \sigma_0^j; \psi \right) \\ &\leq -\frac{1}{N_u} \sum_{j=1}^{N_u} \mathbb{1} \left[ y^j = 0 \right] \log P \left( \sigma_0^j \succ \sigma_1^j; \psi \right) + \mathbb{1} \left[ y^j = 1 \right] \log P \left( \sigma_1^j \succ \sigma_0^j; \psi \right) \\ &\quad + \mathbb{1} \left[ P \left( \sigma_1^j \succ \sigma_0^j; \psi \right) > 0.5, y^j = 0 \right] \log P \left( \sigma_0^j \succ \sigma_1^j; \psi \right) \\ &\quad + \mathbb{1} \left[ P \left( \sigma_0^j \succ \sigma_1^j; \psi \right) > 0.5, y^j = 1 \right] \log P \left( \sigma_1^j \succ \sigma_0^j; \psi \right) \\ &\leq \widehat{\mathcal{L}}'_u(\psi) + \eta_1 \max_j \left\{ -\log P \left( \sigma_0^j \succ \sigma_1^j \right) \right\} + \eta_2 \max_j \left\{ -\log P \left( \sigma_1^j \succ \sigma_0^j \right) \right\} \\ &\leq \widehat{\mathcal{L}}'_u(\psi) + \eta\Omega.\end{aligned}\tag{A.2}$$

**Step 2:** Proving low bound :  $\widehat{\mathcal{L}}_u(\psi) \geq \widehat{\mathcal{L}}'_u(\psi) - \eta\Omega$ .

$$\begin{aligned}\widehat{\mathcal{L}}_u(\psi) &= \frac{1}{N_u} \sum_{j=1}^{N_u} L \left( (\sigma_0^u, \sigma_1^u)^{(j)}, \widehat{y}^j \right) \\ &\geq -\frac{1}{N_u} \sum_{j=1}^{N_u} \mathbb{1} \left[ y^j = 0 \right] \log P \left( \sigma_0^j \succ \sigma_1^j; \psi \right) + \mathbb{1} \left[ y^j = 1 \right] \log P \left( \sigma_1^j \succ \sigma_0^j; \psi \right) \\ &\quad - \mathbb{1} \left[ P \left( \sigma_1^j \succ \sigma_0^j; \psi \right) > 0.5, y^j = 0 \right] \log P \left( \sigma_0^j \succ \sigma_1^j; \psi \right) \\ &\quad - \mathbb{1} \left[ P \left( \sigma_0^j \succ \sigma_1^j; \psi \right) > 0.5, y^j = 1 \right] \log P \left( \sigma_1^j \succ \sigma_0^j; \psi \right) \\ &\geq \widehat{\mathcal{L}}'_u(\psi) - \eta\Omega.\end{aligned}\tag{A.3}$$

756 Combing Step 1 and Step 2, we can obtain the following result:  
757

$$758 \quad \left| \widehat{\mathcal{L}}_u(\psi) - \widehat{\mathcal{L}}'_u(\psi) \right| \leq \eta\Omega. \quad (\text{A.4})$$

759 This completes the proof of Theorem 1.  
760

## 761 A.2 PROOF OF THEOREM 1

762 **Lemma 1.** *Let  $\mathcal{F}$  be a family of functions mapping from  $\mathcal{X}$  to  $\mathbb{R}$  and  $\widehat{\mathcal{D}}$  be empirical datasets  
763 sampled from an i.i.d. sample  $\mathcal{D}$  of size  $N$ . Then, for any  $\delta > 0$ , with probability at least  $1 - \delta$ , the  
764 following holds for all  $f \in \mathcal{F}$ ,*

$$765 \quad \left| \mathbb{E}_{\mathcal{D}}[f] - \mathbb{E}_{\widehat{\mathcal{D}}}[f] \right| \leq 2\widehat{\mathfrak{R}}_{\widehat{\mathcal{D}}}(\mathcal{F}) + 3\sqrt{\frac{\log(2/\delta)}{2N}}, \quad (\text{A.5})$$

766 where  $\mathbb{E}_{\widehat{\mathcal{D}}}[f] = \sum_{i=1}^m [g(x_i)]/m$  is the empirical form of  $\mathbb{E}_{\mathcal{D}}[f]$ . This proof can be found in  
767 Theorem 3.3 of work [Mohri et al. \(2018\)](#).

771 **Proposition 2.** *Let  $\mathcal{F}$  be a family of loss functions defined in Eq. (6) and  $\Pi(\mathcal{R})$  be a family of  
772 functions defined in Definition 2. Then, for any sample  $\widehat{\mathcal{D}} = \{(\sigma_0, \sigma_1, y)^{(i)}\}_{i=1}^N$ , the following  
773 relation holds between the empirical Rademacher complexities of  $\Pi(\mathcal{R})$  and  $\mathcal{F}$ :*

$$774 \quad \widehat{\mathfrak{R}}_{\widehat{\mathcal{D}}}(\mathcal{F}) \leq 2\widehat{\mathfrak{R}}_{\widehat{\mathcal{D}}}(\Pi(\mathcal{R})). \quad (\text{A.6})$$

775 **Proof.** According to Definition 1, the empirical Rademacher complexity of  $\mathcal{F}$  can be written as:  
776

$$\begin{aligned} 777 & \widehat{\mathfrak{R}}_{\widehat{\mathcal{D}}}(\mathcal{F}) \\ 778 & = \mathbb{E}_{\sigma} \left[ \sup_{R \in \mathcal{R}} \frac{1}{N} \sum_{i=1}^N \sigma_i \left[ - (1 - y) \log P(\sigma_0 \succ \sigma_1) - y \log P(\sigma_1 \succ \sigma_0) \right] \right] \\ 779 & \leq \mathbb{E}_{\sigma} \left[ \sup_{R \in \mathcal{R}} \frac{1}{N} \sum_{i=1}^N \sigma_i \left[ - \log P(\sigma_0 \succ \sigma_1) \right] \right] \\ 780 & \quad + \mathbb{E}_{\sigma} \left[ \sup_{R \in \mathcal{R}} \frac{1}{N} \sum_{i=1}^N \sigma_i \left[ - \log P(\sigma_1 \succ \sigma_0) \right] \right] \\ 781 & = 2\widehat{\mathfrak{R}}_{\widehat{\mathcal{D}}}(\Pi(\mathcal{R})). \end{aligned} \quad (\text{A.7})$$

782 This completes the proof of Proposition 2.  
783

791 **The Proof for Theorem 1:** This proof is based on Proposition 1 and 2. Combing Eqs. (A.5) and  
792 (A.6), we can derive that for labeled dataset  $\widehat{\mathcal{D}} = \{(\sigma_0, \sigma_1, y)^{(i)}\}_{i=1}^N$ , the generalization error bound  
793 between expected error  $\mathcal{L}(\psi)$  and empirical error  $\widehat{\mathcal{L}}(\psi)$  holds:  
794

$$795 \quad \mathcal{L}(\psi) \leq \widehat{\mathcal{L}}(\psi) + 4\widehat{\mathfrak{R}}_{\widehat{\mathcal{D}}}(\Pi(\mathcal{R})) + 3\sqrt{\frac{\log(2/\delta)}{2N}}. \quad (\text{A.8})$$

796 However, the dataset for training reward model includes two parts: the labeled dataset  $\mathcal{D}_l =$   
797  $\{(\sigma_0^l, \sigma_1^l, y)^{(i)}\}_{i=1}^{N_l}$  and unlabeled dataset  $\mathcal{D}_u = \{(\sigma_0^u, \sigma_1^u)^{(i)}\}_{i=1}^{N_u}$ , and empirical error  $\widehat{\mathcal{L}}_R(\psi) =$   
798  $\widehat{\mathcal{L}}_l(\psi) + \widehat{\mathcal{L}}_u(\psi)$ . Let  $\widehat{\mathcal{L}}'_u(\psi)$  be the empirical error under unlabeled dataset with true label, then  
799

$$800 \quad \mathcal{L}_R(\psi) \leq \widehat{\mathcal{L}}_l(\psi) + \widehat{\mathcal{L}}'_u(\psi) + 4\widehat{\mathfrak{R}}_{\widehat{\mathcal{D}}}(\Pi(\mathcal{R})) + 3\sqrt{\frac{\log(2/\delta)}{2(N_l + N_u)}}, \quad (\text{A.9})$$

801 where  $\widehat{\mathcal{D}}$  is the input combination of labeled and unlabeled dataset, denoted as  $\{(\sigma_0, \sigma_1)^{(i)}\}_{i=1}^{N_l + N_u}$ .  
802 According to Proposition 1, that is  $\widehat{\mathcal{L}}'_u(\psi) \leq \widehat{\mathcal{L}}_u(\psi) + \eta\Omega$ , we can derive  
803

$$804 \quad \mathcal{L}_R(\psi) \leq \widehat{\mathcal{L}}_R(\psi) + \eta\Omega + 4\widehat{\mathfrak{R}}_{\widehat{\mathcal{D}}}(\Pi(\mathcal{R})) + 3\sqrt{\frac{\log(2/\delta)}{2(N_l + N_u)}}. \quad (\text{A.10})$$

805 This completes the proof of Theorem 1.  
806

## A.3 PROOF OF THEOREM 2

Before proof, we give the related Assumption and Lemmas for Theorem 2. Firstly, we assume the the reward class  $\mathcal{R}$  is realizable (Assumption 1) and bounded (Assumption 3).

**Assumption 3** (Boundedness). *For any  $R \in \mathcal{R}$  and any state-action pairs  $(s, a)$ , the equation  $|R(s, a)| \leq R_{max}$  holds.*

**Proposition 3.** *Consider a set of trajectories  $\{\sigma_0^i, \sigma_1^i\}_{i=1}^N$ , each of length  $L$ , collected from an offline dataset following the distribution  $d_T^\mu(s, a)$ . Then, for any  $\delta \in (0, 1]$ , with probability at least  $1 - \delta$ , the following holds for all  $\hat{R} \in \mathcal{R}$*

$$\left| \mathbb{E}_{(s,a) \sim d_T^\mu(s,a)} [R^*(s, a) - \hat{R}(s, a)] \right| \leq \sqrt{\frac{4C}{NL^2} \log \left( \frac{\mathcal{N}_{\mathcal{R}}(1/N)}{\delta} \right)} + \sqrt{\frac{4R_{max}^2 \log(1/\delta)}{NL}}, \quad (\text{A.11})$$

where  $R^*$  is the true reward model and  $\hat{R}$  is the learned reward model.  $\mathcal{N}_{\mathcal{R}}(1/N)$  is the bracketing number defined in Definition 3.

The proof of Proposition 3 can be found in Appendix A.4. Note that, the above reward model is only trained from offline preference dataset. Through generating more preference data, the upper bound can be tighter. However, the generated data may be inaccuracy, which may improve the reward gap. The selecting mechanism can effectively solve this problem.

**The Proof for Theorem 2:** This Theorem aims to give the lower bound of term  $J(\hat{\pi}, R^*) - J(\mu, R^*)$ . We prove this from ideal case and actual case.

**1) Ideal case:** there is no need to consider distribution shift and empirical error problem. The policy  $\hat{\pi}$  is directly learned by  $\hat{\pi} = \max_{\pi} J(\pi, \hat{R})$ . Thus,  $J(\hat{\pi}, \hat{R}) \geq J(\mu, \hat{R})$  holds. Then, we can derive

$$\begin{aligned} J(\mu, R^*) - J(\hat{\pi}, R^*) &= J(\mu, R^*) - J(\mu, \hat{R}) + J(\mu, \hat{R}) - J(\hat{\pi}, R^*) \\ &\leq J(\mu, R^*) - J(\mu, \hat{R}) + J(\hat{\pi}, \hat{R}) - J(\hat{\pi}, R^*) \\ &\leq |J(\mu, R^*) - J(\mu, \hat{R})| + |J(\hat{\pi}, R^*) - J(\hat{\pi}, \hat{R})|. \end{aligned} \quad (\text{A.12})$$

Since  $J(\pi, R) := \mathbb{E}_{(s,a) \sim d_T^\pi(s,a)} [R(s, a)] / (1 - \gamma)$ ,  $|J(\pi, R^*) - J(\pi, \hat{R})|$  can be written as

$$|J(\pi, R^*) - J(\pi, \hat{R})| = \frac{1}{1 - \gamma} \left| \mathbb{E}_{(s,a) \sim d_T^\pi} [R^*(s, a) - \hat{R}(s, a)] \right|. \quad (\text{A.13})$$

Then, according to Definition 3, we have

$$|J(\pi, R^*) - J(\pi, \hat{R})| \leq \frac{\mathcal{C}_{\mathcal{R}}(\pi)}{1 - \gamma} \left| \mathbb{E}_{(s,a) \sim d_T^\pi} [R^*(s, a) - \hat{R}(s, a)] \right|. \quad (\text{A.14})$$

By Proposition 3, the following inequality holds

$$J(\mu, R^*) - J(\hat{\pi}, R^*) \leq \frac{1 + \mathcal{C}_{\mathcal{R}}(\hat{\pi})}{1 - \gamma} \left( \sqrt{\frac{4C}{NL^2} \log \left( \frac{\mathcal{N}_{\mathcal{R}}(1/N)}{\delta} \right)} + \sqrt{\frac{4R_{max}^2 \log(1/\delta)}{NL}} \right). \quad (\text{A.15})$$

**2) Actual case:** offline RL suffers from distribution shift problem and it is a necessity to consider empirical error. Therefore, many offline RL algorithms incorporate conservatism into policy to overcome distribution shift, that is learning policy through  $\hat{\pi} = \max_{\pi} J(\pi, \hat{R}) - P(\pi)$ . Therefore,  $J(\hat{\pi}, \hat{R}) \geq J(\mu, \hat{R}) - \xi$  instead of  $J(\hat{\pi}, \hat{R}) \geq J(\mu, \hat{R})$  in actual implantation. The term  $\xi$  depends on the algorithm itself. Here, we take CQL as a example. The policy gap is proven in Theorem 3.6 of (Kumar et al., 2020). The detailed is given as follow.

$$\xi = \underbrace{\frac{\gamma C_{T,R} R_{max}}{(1 - \gamma)^2} \mathbb{E}_{s \sim d_T^\mu(s)} \sqrt{\frac{|\mathcal{A}|(1 + D(s))}{|D(s)|}}}_{:=\xi_1} - \underbrace{\frac{\alpha}{1 - \gamma} \mathbb{E}_{s \sim d_T^\mu(s)} [D(s)]}_{:=\xi_2}, \quad (\text{A.16})$$

where  $|\cdot|$  denotes cardinality of a specific set,  $D(s) = \sum_a \hat{\pi}(a|s) \left( \frac{\hat{\pi}(a|s)}{\mu(a|s)} - 1 \right)$  and  $C_{T,R}$  is the empirical coefficient. It consists of two terms: the first term  $\xi_1$  captures the decrease in policy



performance due to sampling error. The second term  $\xi_2$  captures the increase in policy performance in empirical setting (Kumar et al., 2020).

Therefore, when considering distribution shift and empirical error, we can derive

$$\begin{aligned} J(\mu, R^*) - J(\hat{\pi}, R^*) &\leq J(\mu, R^*) - J(\mu, \hat{R}) + J(\hat{\pi}, \hat{R}) + \xi - J(\hat{\pi}, R^*) \\ &\leq \xi + |J(\mu, R^*) - J(\mu, \hat{R})| + |J(\hat{\pi}, R^*) - J(\hat{\pi}, \hat{R})|. \end{aligned} \quad (\text{A.17})$$

Furthermore, according to Eq. (A.15), we have

$$J(\mu, R^*) - J(\hat{\pi}, R^*) \leq \xi + \frac{1 + \mathcal{C}_{\mathcal{R}}(\hat{\pi})}{1 - \gamma} \left( \sqrt{\frac{4C}{NL^2} \log \left( \frac{\mathcal{N}_{\mathcal{R}}(1/N)}{\delta} \right)} + \sqrt{\frac{4R_{max}^2 \log(1/\delta)}{NL}} \right). \quad (\text{A.18})$$

This completes the proof for Theorem 2.

#### A.4 PROOF OF PROPOSITION 3

**Lemma 2.** *There exists an absolute constant  $C$  such that for any  $\delta \in (0, 1]$ , with probability at least  $1 - \delta$ , the following holds for all  $\hat{R} \in \mathcal{R}$*

$$\sum_{i=1}^N \left( P_{\hat{R}}(y^i | \sigma_0^i, \sigma_1^i) - P_{R^*}(y^i | \sigma_0^i, \sigma_1^i) \right)^2 \leq C \log \left( \frac{\mathcal{N}_{\mathcal{R}}(1/N)}{\delta} \right), \quad (\text{A.19})$$

where  $y \in \{0, 1\}$ ,  $R^*$  is the true reward model and  $\hat{R}$  is the learned reward model.  $P_R(0 | \sigma_0, \sigma_1)$  is the probability that  $\sigma_0$  is preferable  $\sigma_1$  under reward model  $R$ .  $\mathcal{N}_{\mathcal{R}}(1/N)$  is the bracketing number defined in Definition 3. This proof can be found in Lemma 2 of previous work (Zhan et al., 2024) and Proposition 14 of (Liu et al., 2022).

**The Proof for Proposition 3:** We prove this Proposition from two steps. Firstly, we bound the difference between  $P_{R^*}(\cdot | \sigma_0, \sigma_1)$  and  $P_{\hat{R}}(\cdot | \sigma_0, \sigma_1)$ . Then, we bound the difference between  $R^*(s, a)$  and  $\hat{R}(s, a)$ .

**Step 1:** Bound the probability difference  $|P_{R^*}(\cdot | \sigma_0, \sigma_1) - P_{\hat{R}}(\cdot | \sigma_0, \sigma_1)|$ .

By Cauchy-Schwarz inequality  $(\sum_i a_i b_i)^2 \leq (\sum_i a_i^2)(\sum_i b_i^2)$ , we set  $a_i = P_{\hat{R}} - P_{R^*}$  and  $b_i$  that is chosen from  $\{-1, 1\}$  and satisfies  $a_i b_i > 0$ . Then, we have

$$\frac{1}{N} \left( \sum_{i=1}^N |P_{\hat{R}}(y^i | \sigma_0^i, \sigma_1^i) - P_{R^*}(y^i | \sigma_0^i, \sigma_1^i)| \right)^2 \leq \sum_{i=1}^N \left( P_{\hat{R}}(y^i | \sigma_0^i, \sigma_1^i) - P_{R^*}(y^i | \sigma_0^i, \sigma_1^i) \right)^2.$$

Then, by Lemma 2, the probability difference can be written as

$$\sum_{i=1}^N |P_{\hat{R}}(y^i | \sigma_0^i, \sigma_1^i) - P_{R^*}(y^i | \sigma_0^i, \sigma_1^i)| \leq \sqrt{CN \log \left( \frac{\mathcal{N}_{\mathcal{R}}(1/N)}{\delta} \right)}. \quad (\text{A.20})$$

**Step 2:** Bound the reward difference  $|R^*(s, a) - \hat{R}(s, a)|$ .

Based on Assumption 3, let  $f(\sigma) = \sum_i R^*(s^i, a^i) \in [-LR_{max}, LR_{max}]$ ,  $g(\sigma) = \sum_i \hat{R}(s^i, a^i) \in [-LR_{max}, LR_{max}]$ , and  $F(x_1, x_2) = e^{x_1} / (e^{x_1} + e^{x_2})$ . Then, according to Eq. (2), we can derive

$$\left| P_{\hat{R}}(y | \sigma_0, \sigma_1) - P_{R^*}(y | \sigma_0, \sigma_1) \right| = \left| F(f(\sigma_0), f(\sigma_1)) - F(g(\sigma_0), g(\sigma_1)) \right|. \quad (\text{A.21})$$

Notably, the above equation only considers the condition  $y = 0$  since the result of  $y = 1$  is similar to  $y = 0$ . Then, according to error propagation, we have

$$\begin{aligned}
& \left| F(f(\sigma_0), f(\sigma_1)) - F(g(\sigma_0), g(\sigma_1)) \right| \\
& \approx \left| \frac{\partial F}{\partial x_1} \right| \left| f(\sigma_0) - g(\sigma_0) \right| + \left| \frac{\partial F}{\partial x_2} \right| \left| f(\sigma_1) - g(\sigma_1) \right| \\
& \geq \frac{e^{x_1+x_2}}{(e^{x_1} + e^{x_2})^2} \left| (f(\sigma_0) - g(\sigma_0)) + (f(\sigma_1) - g(\sigma_1)) \right| \\
& \geq \frac{1}{4} \left| \sum_{l=1}^L \sum_{j=0}^1 (R^*(s_j^l, a_j^l) - \widehat{R}(s_j^l, a_j^l)) \right|.
\end{aligned} \tag{A.22}$$

Then, combing Eqs. (A.20) and (A.22), we can further derive

$$\left| \sum_{i=1}^N \sum_{l=1}^L \sum_{j=0}^1 (R^*(s_j^{i,l}, a_j^{i,l}) - \widehat{R}(s_j^{i,l}, a_j^{i,l})) \right| \leq 4 \sqrt{CN \log \left( \frac{\mathcal{N}_{\mathcal{R}}(1/N)}{\delta} \right)}. \tag{A.23}$$

According to Chernoff-Hoeffding bound (Hoeffding, 1994), for independent random variables  $X_1, X_2, \dots, X_n$ , with high probability  $1 - \delta$ , the below equation holds

$$\mathbb{E}[X] \leq \frac{1}{n} \sum_{i=1}^n X_i + (b - a) \sqrt{\frac{\log(1/\delta)}{2n}}, \tag{A.24}$$

where  $[a, b]$  is the range of values that each  $X_i$  can take. Then, for Eq. (equation A.23), we set  $X_i$  as  $R^*(s_j, a_j) - R(s_j, a_j)$  and  $X_i \in [-2R_{max}, 2R_{max}]$ . Then, the below equation holds

$$\left| \mathbb{E}_{(s,a) \sim d_T^{\mu}(s,a)} [R^*(s, a) - \widehat{R}(s, a)] \right| \leq \sqrt{\frac{4C}{NL^2} \log \left( \frac{\mathcal{N}_{\mathcal{R}}(1/N)}{\delta} \right)} + \sqrt{\frac{4R_{max}^2 \log(1/\delta)}{NL}}, \tag{A.25}$$

where  $d_T^{\mu}(s, a)$  is the distribution of offline dataset. The above equation holds since the trajectories  $\{\sigma_0^i, \sigma_1^i\}_{i=1}^N$  are collected from offline dataset. This completes the proof of Proposition 3.

## B RELATED EXPERIMENTS

We conduct experiments on Mujoco and Adroit tasks, which are included in the D4RL (Fu et al., 2020) benchmark. The code for LEASE is available at [github.com/\\*\\*\\*\\*](https://github.com/****). This part introduces detailed experiments setup, hyper-parameter and parameter analysis for LEASE.

### B.1 EXPERIMENTS SETUP

**Offline dataset.** We employ Mujoco and Adroit environments to test the performance of LEASE. The Mujoco tasks include halfcheetah-v2, hopper-v2, and walker2d-v2. The Adroit tasks include pen-v1, door-v1 and hammer-v1. Fig. A.1 depicts the above six tasks. For three Mujoco tasks, the goal is to control the robot’s various joints to achieve faster and more stable locomotion. For Adroit tasks, they involve controlling a 24-DoF simulated Shadow Hand robot on a sparse reward, high-dimensional robotic manipulation task (Fu et al., 2020). The offline datasets are based on the D4RL dataset. We select medium and medium-expert two types dataset for Mujoco tasks and human and expert for Adroit tasks. The difference between different dataset in certain task lies in the collected policy.

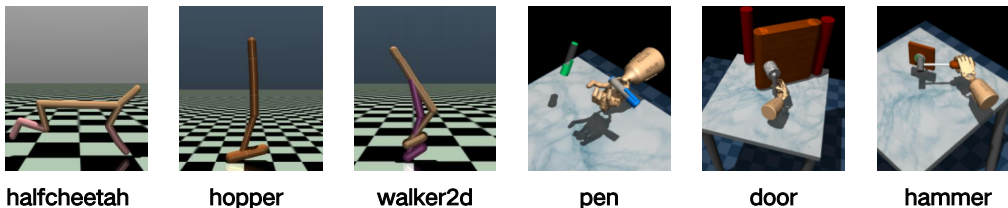


Figure A.1: The description of Mujoco tasks (version 2) and Adroit tasks (version 1).

**Preference dataset.** The preference dataset is from (Yuan et al., 2024). We use two types preference dataset to test algorithm performance. One is that the labels for preference are from human feedback, the another is the labels are from the ground-truth reward. Obviously, since human preferences are subjective, the labels from ground-truth reward do not correspond to human preference in actual situation, and the performance of preference labels calibrated by ground-truth reward is better than that of labels calibrated by human feedback in most cases. Here, we use the labels from human feedback in Mujoco tasks and from ground-truth reward in Adroit tasks.

**Basic offline algorithms.** Offline PbRL setting involves two steps: reward learning and policy learning. Here, we directly use offline RL algorithms CQL (Kumar et al., 2020) and IQL (Kostrikov et al., 2021b) to perform policy learning. The performance of these algorithms is evaluated based on the cumulative return. For comparison, the scores are normalized between 0 (random policy score) and 100 (expert policy score) (Fu et al., 2020). The normalized score  $\tilde{S}$  is computed by:

$$\tilde{S} = \frac{S - S_r}{S_e - S_r} \times 100,$$

where  $S_r$ ,  $S_e$  and  $S$  are the expected return of a random policy, an expert policy, and the trained policy by offline RL algorithms, respectively.

## B.2 HYPER-PARAMETERS

**Reward and transition model training.** Table A.1 gives the hyperparameter configuration for reward and transition training. Similar to (Yu et al., 2020), we train an ensemble of 7 transition models and select the best 5 models. Each model consists of a 4-layer feed-forward neural network with 200 hidden units. The model training employs maximum likelihood estimation with a learning rate of  $1e^{-3}$  and Adam optimizer. Following (Yuan et al., 2024), we train an ensemble of 3 reward models. Each reward model includes a 3-layer feed-forward neural network with 256 hidden units. The reward model training uses cross-entropy loss with a learning rate of  $3e^{-4}$  and Adam optimizer.

Table A.1: Base parameter configuration of LEASE.

| Model            | Parameter                        | Value             |
|------------------|----------------------------------|-------------------|
| Transition Model | Model learning rate              | $1 \times e^{-3}$ |
|                  | Number of hidden layers          | 4                 |
|                  | Number of hidden units per layer | 200               |
|                  | Batch size                       | 256               |
|                  | Number of model networks $N_T$   | 7                 |
| Reward Model     | Number of elites                 | 5                 |
|                  | Model learning rate              | $3 \times e^{-4}$ |
|                  | Number of hidden layers          | 3                 |
|                  | Number of hidden units per layer | 256               |
|                  | Batch size (Pretrain)            | 256(64)           |
|                  | Number of model networks $N_R$   | 3                 |
|                  | Number of labeled dataset $N_l$  | 100               |

Table A.2: The hyperparameters for LEASE under CQL and IQL algorithms.

| Task Name                 | CQL |            |               | IQL |            |               |
|---------------------------|-----|------------|---------------|-----|------------|---------------|
|                           | $H$ | $\kappa_p$ | $\kappa_\tau$ | $H$ | $\kappa_p$ | $\kappa_\tau$ |
| walker2d-medium           | 10  | 0.85       | 0.05          | 10  | 0.75       | 0.08          |
| walker2d-medium-expert    | 10  | 0.85       | 0.05          | 10  | 0.85       | 0.08          |
| hopper-medium             | 10  | 0.85       | 0.05          | 10  | 0.85       | 0.08          |
| hopper-medium-expert      | 100 | 0.85       | 0.05          | 10  | 0.85       | 0.08          |
| halfcheetah-medium        | 10  | 0.85       | 0.05          | 10  | 0.85       | 0.08          |
| halfcheetah-medium-expert | 200 | 0.85       | 0.05          | 10  | 0.75       | 0.08          |
| pen-human                 | 10  | 0.95       | 0.05          | 10  | 0.85       | 0.08          |
| pen-expert                | 10  | 0.90       | 0.05          | 10  | 0.95       | 0.08          |
| door-human                | 10  | 0.95       | 0.05          | 10  | 0.90       | 0.08          |
| door-expert               | 20  | 0.95       | 0.05          | 10  | 0.95       | 0.08          |
| hammer-human              | 10  | 0.95       | 0.05          | 10  | 0.90       | 0.08          |
| hammer-expert             | 20  | 0.99       | 0.05          | 10  | 0.95       | 0.08          |

**Policy optimization.** Offline PbRL involves reward learning and policy optimization. For LEASE, the reward is updated during policy learning. Therefore, the reward is influenced by the rollout length and rollout batch size and current learned policy. Here, we select CQL and IQL as basic offline algorithms. In CQL algorithm, `cql_weight` and `temperature` are set as 5.0 and 1.0 for all tasks, respectively. In IQL algorithm, `expectile` and `temperature` are set as 0.7 and 3.0 for mujoco tasks and 0.8 and 3.0 for adroit tasks.

The critic network  $Q_\omega$  and the policy network  $\pi_\theta$  adopt a 2-layer feed-forward neural network with 256 hidden units. The hyperparameters for LEASE includes: rollout horizon  $H \in \{10, 20, 100, 200\}$ , probability confidence  $\kappa_p(\sigma_0^u, \sigma_1^u, \hat{y}) \in \{0.75, 0.85, 0.90, 0.95, 0.99\}$  and uncertainty variance  $\kappa_\tau(\sigma_0^u, \sigma_1^u, N_R) \in \{0.05, 0.08\}$ . Table A.2 shows the above hyperparameters for LEASE under CQL and IQL algorithms. The maximum buffer capacity of unlabeled data is 50000. In next part , we give detailed analysis for the above hyper-parameters.

### B.3 PARAMETER ANALYSIS

**Number of preference data  $N_l$ .** The number of preference data  $N_l$  influences the accuracy of reward model directly. URLHF (Yuan et al., 2024) trains reward model with 2000 preference dataset. Our method LEASE aims to achieve comparable performance with URLHF under fewer preference dataset. To test the effects of the number of labeled preference data for agent performance, we analyze the performance when the number of preference data  $N_l$  is 20. The hyper-parameters of 20 preference dataset are same with that of 100 preference dataset, apart from the number of preference dataset and the batch size of pretrained reward model. The batch size is set as 16 here.

Table A.3: The comparison results for the D4RL tasks under different number of preference data  $N_l$ .  $\pm$  captures the standard deviation over seeds. Bold indicates the highest score.

| Task Name                 | CQL (Kumar et al., 2020)         |                                   | IQL (Kostrikov et al., 2021b) |                                   |                                   |
|---------------------------|----------------------------------|-----------------------------------|-------------------------------|-----------------------------------|-----------------------------------|
|                           | $N_l = 20$                       | $N_l = 100$                       | $N_l = 20$                    | $N_l = 100$                       |                                   |
| walker2d-medium           | 77.7 $\pm$ 1.3                   | <b>78.4 <math>\pm</math> 0.9</b>  | pen-human                     | 71.9 $\pm$ 7.4                    | <b>75.6 <math>\pm</math> 3.3</b>  |
| walker2d-medium-expert    | 98.0 $\pm$ 18.6                  | <b>98.6 <math>\pm</math> 18.1</b> | pen-expert                    | 102.5 $\pm$ 12.8                  | <b>113.8 <math>\pm</math> 6.3</b> |
| hopper-medium             | <b>56.8 <math>\pm</math> 1.5</b> | 56.5 $\pm$ 0.6                    | door-human                    | 4.4 $\pm$ 1.2                     | <b>5.9 <math>\pm</math> 0.5</b>   |
| hopper-medium-expert      | 54.5 $\pm$ 1.2                   | <b>56.4 <math>\pm</math> 0.8</b>  | door-expert                   | <b>105.2 <math>\pm</math> 0.1</b> | <b>105.2 <math>\pm</math> 0.2</b> |
| halfcheetah-medium        | 43.4 $\pm$ 0.4                   | <b>43.5 <math>\pm</math> 0.1</b>  | hammer-human                  | 1.2 $\pm$ 0.4                     | <b>1.7 <math>\pm</math> 0.4</b>   |
| halfcheetah-medium-expert | 51.0 $\pm$ 0.8                   | <b>53.2 <math>\pm</math> 3.1</b>  | hammer-expert                 | <b>126.4 <math>\pm</math> 0.1</b> | 126.3 $\pm$ 0.1                   |
| <b>Mujoco Average</b>     | 63.6 $\pm$ 4.0                   | <b>64.4 <math>\pm</math> 4.0</b>  | <b>Adroit Average</b>         | 68.6 $\pm$ 3.7                    | <b>71.4 <math>\pm</math> 1.8</b>  |

Table A.3 shows the agent performance when the number of preference dataset is 20 and 100, where we use CQL algorithm for mujoco tasks and IQL algorithm for adroit tasks. This table indicates LEASE still can improve agent performance under very fewer preference dataset. However, as the number of preference data  $N_l$  decreases, the average of agent performance is slightly reduced. This is mainly because the reward model pretrained with very little preference data has a large generalization error, resulting in low quality of the generated preference data, which in turn affects

the accuracy of the final reward model. Notably, the performance gap between 100 preference data and 20 preference data is not much. This indicates that LEASE have potential to perform well under more fewer preference dataset.

**Rollout length  $H$ .** The rollout length  $H$  is equal to the length of generated unlabeled data  $\widehat{L}$ . The  $H$  is influenced by the accuracy of trained dynamics model. As  $H$  increases, the accuracy of prediction decreases and large error will bring the unstable agent training. However, the long horizon of trajectory can better describe the preference of human feedback and can reduce sample complexity (Eq. (16)). Table A.4 gives the prediction accuracy of the trained reward model for different rollout length  $H$ . Here, we take CQL as example and select the last 200 of 2000 preference data as evaluation dataset that are not seen in in reward training stage.

The accuracy is calculated through the gap between predicted preference and true preference based on two trajectories. It shows that long horizon easily brings low performance of reward model since the cumulative error of the trained transition model, but long horizon is beneficial to enhance the performance of reward model under some conditions. The accuracy of preference model can be improved through data augmentation under most tasks. Note that the result can not completely represent the generalization performance of reward model since the evaluation data are limited. The choice of rollout length  $H$  mainly depends on the accuracy of the trained transition model and the agent performance.

Table A.4: The prediction result of preference model for different rollout length  $H$ . The preference model is based on reward model. The performance of preference model can directly reflect the performance of reward model. The evaluation data are not seen in reward training stage.

| Task name               | walker2d-m  | walker2d-m-e | hopper-m    | hopper-m-e  | halfcheetah-m | halfcheetah-m-e |
|-------------------------|-------------|--------------|-------------|-------------|---------------|-----------------|
| <b>Pretrained model</b> | <b>0.59</b> | 0.85         | 0.73        | 0.77        | 0.6           | 0.72            |
| <b>Updated</b> $H = 10$ | 0.58        | <b>0.85</b>  | <b>0.78</b> | 0.74        | <b>0.62</b>   | <b>0.78</b>     |
| $H = 100$               | 0.57        | 0.84         | 0.71        | <b>0.79</b> | 0.57          | 0.74            |
| $H = 200$               | 0.57        | 0.83         | 0.74        | 0.76        | 0.60          | 0.77            |

**Probability confidence  $p$  and uncertainty variance  $\tau$ .** The above two parameters influence the selection of unlabeled preference dataset (Eq. (8)). For reward model composed of simple fully connected layers, using excessive amounts of unlabeled data for training can easily lead to over-fitting. Conversely, using too little data can result in poor generalization of the model. Implementing a selection mechanism ensures data quality on one hand while preventing the reward model from over-fitting on the other. Since the label is from the ground-truth reward for adroit tasks, the pre-trained reward model is closed to true model. Therefore, the probability confidence  $p$  in adroit tasks is set higher than that of Mujoco tasks.

#### B.4 ADDITIONAL RESULTS

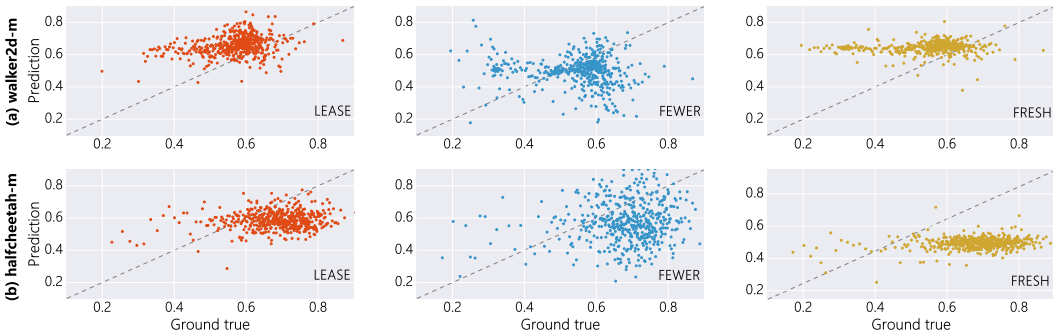
**The other comparison results for screen mechanism.** Table A.5 shows the performance improvement under Adroit tasks when using screen mechanism  $f(\sigma_0, \sigma_1)$ . It shows that the performance can be significantly improved under CQL algorithm through selecting mechanism, but not significant under IQL algorithm. In general, combing with Table 2, we can conclude the selecting mechanism for unlabeled dataset can effectively improve the agent performance.

Table A.5: The comparison results between the performance using selecting mechanism and that not using. The latter method is denoted as FRESH.  $\uparrow$  denotes the improvement of performance.

| Method       | CQL (Kumar et al., 2020)        |                                |                                 | IQL (Kostrikov et al., 2021b) |                         |                         |
|--------------|---------------------------------|--------------------------------|---------------------------------|-------------------------------|-------------------------|-------------------------|
|              | pen-expert                      | door-expert                    | hammer-expert                   | pen-expert                    | door-expert             | hammer-expert           |
| <b>FRESH</b> | 113.2 $\pm$ 14.8                | 101.6 $\pm$ 3.5                | 99.5 $\pm$ 34.5                 | 113.4 $\pm$ 14.7              | 105.1 $\pm$ 0.1         | 125.5 $\pm$ 0.8         |
| <b>LEASE</b> | <b>132.5</b> $\uparrow$ (17.0%) | <b>103.2</b> $\uparrow$ (1.6%) | <b>126.3</b> $\uparrow$ (26.5%) | 113.8 $\uparrow$ (0.3%)       | 105.2 $\uparrow$ (0.2%) | 126.3 $\uparrow$ (0.7%) |

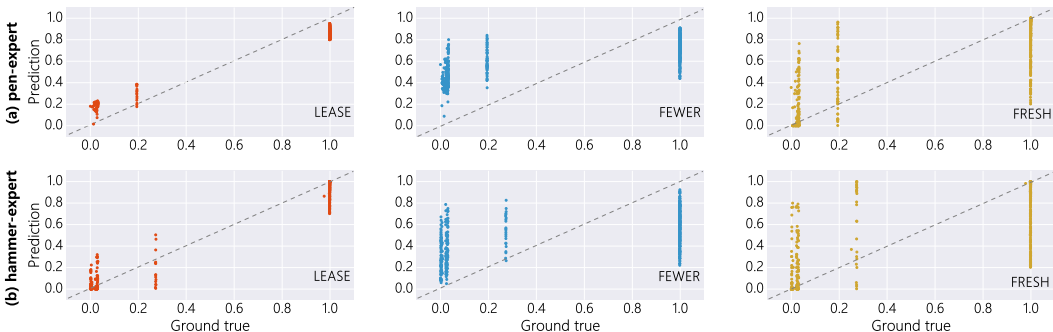
**The other comparison results for reward model performance.** Fig. A.2 and A.3 show comparison between prediction value by the learned rewards and their ground truths for LEASE, FEWER and FRESH under other Mujoco and Adroit tasks, where the offline algorithm is IQL. The predicted and true rewards are both normalized to  $[0, 1]$ . We randomly sample 500 data from unlabeled datasets that are not seen in training stage for evaluation like Fig. 3. Since the reward has the value in

1134  
1135  
1136  
1137  
1138  
1139  
1140  
1141  
1142  
1143  
1144  
1145



1146 Figure A.2: The comparison between prediction value by the learned rewards and their ground truths  
1147 for different methods under (a) walker2d-medium and (b) halfcheetah-medium datasets.

1149  
1150  
1151  
1152  
1153  
1154  
1155  
1156  
1157  
1158  
1159  
1160



1162 Figure A.3: The comparison between prediction value by the learned rewards and their ground truths  
1163 for different methods under (a) pen-expert and (b) hammer-expert datasets.

1164  
1165  
1166  
1167  
1168  
1169  
1170  
1171  
1172  
1173

certain position in Adroit tasks, the figure is in the form of multiple scattered lines in Fig. A.3. This also shows that the data augmentation and selecting mechanism can improve the reward model performance.

**Comparison results between LEASE and URLHF with fewer preference dataset.** To further show superior performance of the proposed method, Table A.6 compares LEASE to the baseline algorithm URLHF with the same amount of data as LEASE, where the latter method is denoted as URLHF\*. This table shows that the average performance of LEASE is superior to that of the baseline algorithm URLHF using the same amount data with LEASE.

1174  
1175  
1176  
1177  
1178  
1179  
1180  
1181  
1182  
1183  
1184  
1185  
1186  
1187

Table A.6: Comparison results between URLHF using fewer preference data and LEASE, where the offline RL algorithm is CQL and URLHF using fewer data is denoted as URLHF\*.

| Task Name       | URLHF*             | URLHF              | FEWER             | LEASE              |
|-----------------|--------------------|--------------------|-------------------|--------------------|
| walker2d-m      | 76.1 ± 0.8         | 76.0 ± 0.9         | 77.4 ± 0.6        | <b>78.4 ± 0.9</b>  |
| walker2d-m-e    | 86.8 ± 18.6        | 92.8 ± 22.4        | 77.7 ± 0.3        | <b>98.6 ± 18.1</b> |
| hopper-m        | <b>56.6 ± 2.4</b>  | 54.7 ± 3.4         | 55.8 ± 2.8        | <b>56.5 ± 0.6</b>  |
| hopper-m-e      | 55.3 ± 0.9         | <b>57.4 ± 4.9</b>  | 53.6 ± 0.9        | 56.4 ± 0.8         |
| halfcheetah-m   | 43.3 ± 0.2         | <b>43.4 ± 0.1</b>  | <b>43.5 ± 0.1</b> | <b>43.5 ± 0.1</b>  |
| halfcheetah-m-e | 58.9 ± 2.3         | <b>62.7 ± 7.1</b>  | 48.3 ± 0.7        | 53.2 ± 3.1         |
| Mujoco Average  | 62.8 ± 4.2         | <b>64.5 ± 6.5</b>  | 59.4 ± 0.9        | <b>64.4 ± 4.0</b>  |
| pen-human       | <b>17.7 ± 13.0</b> | 9.8 ± 14.1         | 0.5 ± 3.0         | 3.8 ± 4.6          |
| pen-expert      | 114.6 ± 53.7       | <b>138.3 ± 5.2</b> | 128.1 ± 0.7       | 132.5 ± 2.3        |
| door-human      | 1.7 ± 1.1          | <b>4.7 ± 5.9</b>   | 0.2 ± 1.0         | <b>4.7 ± 8.8</b>   |
| door-expert     | <b>103.3 ± 0.5</b> | <b>103.9 ± 0.8</b> | 103.0 ± 0.9       | <b>103.2 ± 0.7</b> |
| hammer-human    | 0.7 ± 0.1          | <b>0.9 ± 0.3</b>   | 0.3 ± 0.0         | 0.3 ± 0.0          |
| hammer-expert   | 117.4 ± 2.7        | 120.2 ± 6.8        | 124.1 ± 2.1       | <b>126.3 ± 1.2</b> |
| Adroit Average  | 59.2 ± 11.8        | <b>63.0 ± 5.5</b>  | 59.4 ± 1.3        | 61.8 ± 3.0         |

**The effect of the introduction of uncertainty for pseudo-labeling accuracy.** Table A.7 validates the advantage of using uncertainty for reducing pseudo-labeling error, where the accuracy of pseudo-label generated by reward model is evaluated on all preference dataset. The below table shows that using uncertainty can improve accuracy of pseudo labels.

Table A.7: Comparison results of pseudo-labeling accuracy between using confidence and uncertainty and only using confidence.

| Task name                         | pen-expert | door-expert | hammer-expert |
|-----------------------------------|------------|-------------|---------------|
| <b>confidence and uncertainty</b> | 87.25%     | 89.25%      | 85.45%        |
| <b>only confidence</b>            | 85.85%     | 87.80%      | 84.41%        |

**The other results for LEASE under model-based offline algorithm.** Table A.8 shows the results of COMBO (Yu et al., 2021) under designed framework. For COMBO hyperparameters, the rollout horizon of preference trajectory  $H$ , probability confidence  $\kappa_p$  and uncertainty variance  $\kappa_\tau$  are set as 10, 0.85 and 0.05 for all tasks, respectively. Note that in our framework, model-based methods do not necessarily perform better than model-free methods. Model-based RL methods focus on how to learn conservative policy by regularizing  $Q$  values or penalizing rewards to alleviate the effects of inaccuracy model data. Therefore, model-based RL requires higher accuracy of the reward model than model-free RL.

Table A.8: Comparison results of offline RL algorithms under the designed framework.

| Task Name             | CQL*            | IQL*            | COMBO*         |
|-----------------------|-----------------|-----------------|----------------|
| walker2d-m            | 78.4 $\pm$ 0.9  | 74.6 $\pm$ 1.8  | 71.6 $\pm$ 2.4 |
| walker2d-m-e          | 98.6 $\pm$ 18.1 | 108.1 $\pm$ 0.5 | 79.1 $\pm$ 1.1 |
| hopper-m              | 56.5 $\pm$ 0.6  | 56.0 $\pm$ 0.5  | 54.8 $\pm$ 0.9 |
| hopper-m-e            | 56.4 $\pm$ 0.8  | 55.9 $\pm$ 1.9  | 54.9 $\pm$ 1.1 |
| halfcheetah-m         | 43.5 $\pm$ 0.1  | 43.0 $\pm$ 0.3  | 42.9 $\pm$ 0.1 |
| halfcheetah-m-e       | 53.2 $\pm$ 3.1  | 62.4 $\pm$ 1.4  | 73.8 $\pm$ 7.0 |
| <b>Mujoco Average</b> | 64.4 $\pm$ 4.0  | 66.7 $\pm$ 1.0  | 62.9 $\pm$ 2.1 |

**The analysis of the accuracy of transition model for agent performance.** Table A.9 shows the detailed results for the effect of transition model accuracy for agent performance. It shows that the lower the accuracy of the transition model, the poorer performance of the agent, where the accuracy of transition model is measured by the sum of the mean square values of the predicted value and the true value in each dimension.

Table A.9: Results of the effect of the learned transition model for agent performance.

| hopper-medium          |                   | pen-expert             |                   |
|------------------------|-------------------|------------------------|-------------------|
| Transition model error | Agent performance | Transition model error | Agent performance |
| 0.35 $\pm$ 0.01        | 56.5 $\pm$ 0.60   | 1.05 $\pm$ 0.01        | 132.5 $\pm$ 2.3   |
| 0.49 $\pm$ 0.04        | 54.8 $\pm$ 0.85   | 1.42 $\pm$ 0.05        | 126.4 $\pm$ 8.63  |
| 1.19 $\pm$ 0.05        | 52.85 $\pm$ 0.92  | 2.36 $\pm$ 0.24        | 87.65 $\pm$ 41.79 |

## C FURTHER DISCUSSION

### C.1 DISCUSSION FOR SAMPLE EFFICIENCY

**The discussion of sample efficiency.** In offline RL field, high sample efficiency refers that the agent can achieve comparable performance under fewer data compared with the performance under large data. In this paper, the data refers to the preference dataset. The labeled preference dataset, each trajectory of length  $L$ , is collected through real-time human feedback under policy  $\mu$ , which demands tremendous human effort, thus the collected cost of preference data is higher than fixed offline data. The unlabeled dataset is generated through trained transition without real-time interaction under learned policy  $\pi^t$  at time  $t$ . Through data augmentation, the sample efficiency can be significantly reduced. The performance gap caused by reward can be reduced to  $\varpi$  under fewer labeled dataset.

**Comparison with (Zhan et al., 2024).** Zhan et al. (2024) developed systematical theory for offline PbRL and also introduced the concentrability coefficient for PbRL, defined as

$$\mathcal{C}_r(\pi) = \max \left\{ 0, \sup_r \frac{\mathbb{E}_{\sigma^0 \sim \pi, \sigma^1 \sim \pi_{\text{ref}}} [r^*(\sigma^0, \sigma^1) - r(\sigma^0, \sigma^1)]}{\sqrt{\mathbb{E}_{\sigma^0 \sim \mu_0, \sigma^1 \sim \mu_1} |r^*(\sigma^0, \sigma^1) - r(\sigma^0, \sigma^1)|^2}} \right\}. \quad (\text{A.26})$$

where  $r(\sigma^0, \sigma^1) = \sum_i [R(s_i^0, a_i^0) - R(s_i^1, a_i^1)]$ ,  $\pi_{\text{ref}}$  is an arbitrary trajectory distribution (usually set as  $\mu_1$ ), and  $\mu_0, \mu_1$  are behavior trajectory distribution. However, it is based on trajectory and is difficult to combine with other offline RL theories. The concentrability coefficient defined in Eq. (14) is based on state-action pairs.

Moreover, the performance gap of offline PbRL between behavior policy and learned policy is influenced by offline algorithm itself and the performance of the learned reward (preference) model. However, Zhan et al. (2024) fails to consider the gap caused by offline algorithm itself. The theory developed in our paper can be easily combined with other offline algorithm. The method in (Zhan et al., 2024) can learn  $\varpi$ -optimal policy with a sample complexity of

$$N = \tilde{O} \left( \frac{c^2 k^2 \mathcal{C}_r^2(\hat{\pi})}{\varpi^2} \log \left( \frac{\mathcal{N}_r(1/N)}{\delta} \right) \right), \quad (\text{A.27})$$

where  $c > 0$  is a universal constant,  $k = (\inf_{x \in [-r_{max}, r_{max}]} \Phi'(x))^{-1}$ , and  $\Phi(x)$  is a monotonically increasing link function. Compared with Eq. (A.27), the sample complexity of LEASE contains more useful information. For example, the sample complexity can be reduced when the length of preference data or the learned policy is closed to behavior policy.

## C.2 DISCUSSION FOR RELATED WORKS

**Model-free offline RL.** Existing model-free algorithms typically use two approaches: policy constraint and value regularization. The goal of policy constraint is to keep the learned policy close to the behavior policy (Kumar et al., 2019). On the other hand, value regularization methods mitigate the value overestimation for out-of-distribution (OOD) data by conservatively estimating  $Q$  values (Kumar et al., 2020) or by penalizing based on the uncertainty of the  $Q$  function (An et al., 2021). However, the core challenge for offline RL arises from limited data coverage. Model-free offline RL algorithms can only learn policies from the offline dataset, which restricts the agent’s ability to explore.

**Model-based offline RL.** Model-based offline RL algorithms train a dynamics model using the offline dataset and leverage this model to enhance data coverage. However, due to the limitations of the offline dataset, there is a discrepancy between the learned dynamics model and the actual dynamics. To address this, conservatism should be integrated into the algorithms to prevent the agent from operating in areas where the predictions of the learned dynamics model are unreliable. One approach is to penalize the reward based on uncertainty quantification (Yu et al., 2020; Sun et al., 2023). Another approach is to enforce lower  $Q$ -values for data generated by the dynamics model that are deemed imprecise (Yu et al., 2021; Liu et al., 2023).

**Comparison with Surf (Park et al., 2022).** Surf is the PbRL method similar to LEASE using data augmentation technique. The differences between them mainly includes the below three aspects: 1) Surf belongs to online RL, but LEASE belongs to offline RL. Surf generates data through interaction with environment (simulator) while LEASE generates data through the learned transition model; 2) Surf only uses confidence for label filtering, whereas LEASE employs both confidence and uncertainty principles, which effectively reduce pseudo-label error; 3) LEASE provides the general theoretical framework for offline PbRL, but Surf don’t provide theoretical analysis.

## C.3 BROADER IMPACTS

**Actual application.** In actual scenarios, especially for human-in-loop-control, such as exoskeleton robot assistance or rehabilitation, designing a high sample efficient RL algorithm is greatly significant. Firstly, many rewards are difficult to describe in mathematical terms in some situation, such as the comfort of interaction. Human feedback or preference is the better way to reflect the above indexes. In addition, in human in-the-loop control, obtaining preference data through interaction



1296 can easily cause fatigue, and inappropriate interaction may cause damage to human. It is a necessity  
1297 to learn reward model from limited preference dataset. Therefore, PbRL has a great potential to  
1298 improve control performance in some scenarios where the reward function is difficult to describe  
1299 and human is in the control loop.

1300 **Theoretical study.** The theory of LEASE shows that the performance gap between the behavior  
1301 policy  $\mu$  and  $\hat{\pi}$  learned by LEASE includes two part: the gap  $\xi$  caused by offline algorithm itself  
1302 and the gap  $\xi_1$  caused by reward model gap (the details see Theorem 2), where the term  $\xi_1$  only  
1303 depends on preference dataset. Moreover, the theory of LEASE is based on state-action pairs, which  
1304 is consistent with most offline algorithms theory. Therefore, the theory developed in this paper can  
1305 be easily combined with other offline algorithms theory and easily used to build theory of policy  
1306 improvement guarantee, which can provide the theoretical basis for offline PbRL and facilitate the  
1307 further development of offline PbRL theory.

1308

1309

1310

1311

1312

1313

1314

1315

1316

1317

1318

1319

1320

1321

1322

1323

1324

1325

1326

1327

1328

1329

1330

1331

1332

1333

1334

1335

1336

1337

1338

1339

1340

1341

1342

1343

1344

1345

1346

1347

1348

1349