

# AUTHORS

Aylin Rasteh  
Pouya Lahabi

# ADVANCEMENTS IN OFFLINE REINFORCEMENT LEARNING



Adversarial Model for Offline Reinforcement Learning  
Model-based Offline Policy Optimization  
Conservative Offline Model-Based Policy Optimization

## INTRODUCTION

Offline Reinforcement Learning (RL) uses pre-collected datasets to train policies without new interactions with the environment, which is crucial for applications where gathering data is costly or risky. This approach is vital in fields like autonomous driving and healthcare, where safety and reliability are essential. Offline RL enables robust policy deployment in real-world scenarios but faces challenges like distributional shift, uncertainty estimation, and ensuring policy robustness.

## OBJECTIVE

The objective of this poster is to present and compare the methodologies, theoretical insights, and empirical results from three key papers in offline reinforcement learning: Adversarial Model for Offline Reinforcement Learning (ARMOR), Model-based Offline Policy Optimization (MOPO), and Conservative Offline Model-Based Policy Optimization (COMBO)

## RELATED LITERATURE

1. Bhardwaj, M., Xie, T., Boots, B., Jiang, N., & Cheng, C.-A. (2023). Adversarial Model for Offline Reinforcement Learning.
2. Yu, T., Thomas, G., Yu, L., Ermon, S., Zou, J., Levine, S., Finn, C., & Ma, T. (2020). Model-based Offline Policy Optimization.
3. Yu, T., Kumar, A., Rafailov, R., Rajeswaran, A., Levine, S., & Finn, C. (2020). Conservative Offline Model-Based Policy Optimization.

## METHODOLOGY

### Adversarial Model for Offline Reinforcement Learning (ARMOR)

- Adversarial training for robust policy improvement.
- Ensures performance relative to a reference policy under uncertain data conditions.
- Achieves state-of-the-art results with a single model.

### Model-based Offline Policy Optimization (MOPO)

- Incorporates model uncertainty into the reward function.
- Penalizes rewards based on dynamics model uncertainty.
- Provides a theoretical guarantee for a lower bound on policy performance.
- Outperforms traditional model-free and other model-based offline RL algorithms.

### Conservative Offline Model-Based Policy Optimization (COMBO)

- Integrates conservative planning into the learning process.
- Enhances robustness and reliability of policy optimization.
- Achieves superior performance in complex dynamic environments.

## ARMOR

We propose to directly optimize **relative performance** to the reference policy

$$\max_{\pi \in \Pi} J(\pi) - J(\pi_{\text{ref}})$$

to obtain a lower bound that is tight at  $\pi_{\text{ref}}$

ARMOR constructs a two-player game using hypothesis class  $\mathcal{M}_\alpha$

$$\hat{\pi} \in \max_{\pi \in \Pi} \left[ \min_{M \in \mathcal{M}_\alpha} J_M(\pi) - J_M(\pi_{\text{ref}}) \right]$$

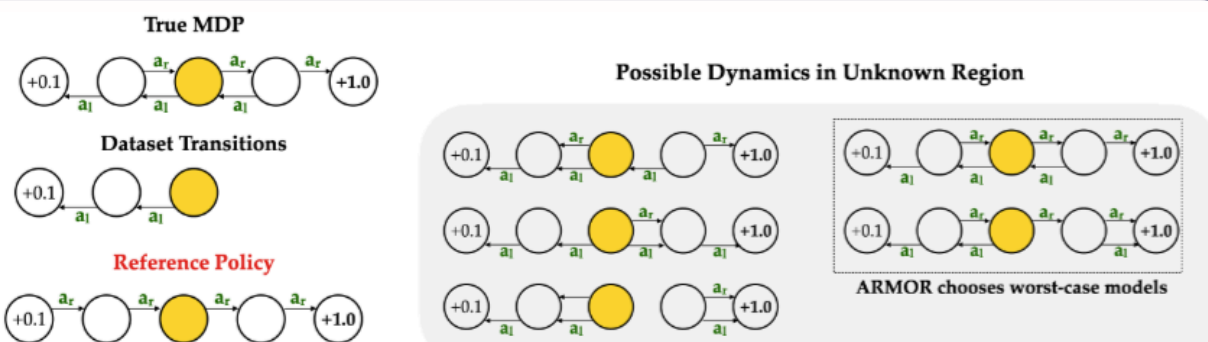
where inner optimization minimizes lower bound on relative performance if  $\mathcal{M}_\alpha$  contains true MDP model  $M^*$

$$\mathcal{M}_\alpha := \{M \in \mathcal{M} : \mathcal{E}_D(M) - \min_{M' \in \mathcal{M}} \mathcal{E}_D(M') \leq \alpha\}$$

$$\mathcal{E}_D(M) := \sum_D -\log P_M(s'|s, a) + (R_M(s, a) - r)^2 / V_{\max}^2$$

Choose  $\alpha$  s.t.  $M^* \in \mathcal{M}_\alpha$

ARMOR ensures policy improvement even when  $\pi_{\text{ref}}$  is not covered by the dataset



### Algorithm 1 ARMOR (Adversarial Model for Offline Reinforcement Learning)

**Input:** Batch data  $\mathcal{D}_{\text{real}}$ , policy  $\pi$ , MDP model  $M$ , critics  $f_1, f_2$ , horizon  $H$ , constants  $\beta, \lambda \geq 0$ ,  $\tau \in [0, 1]$ ,  $w \in [0, 1]$ .

- 1: Initialize target networks  $\bar{f}_1 \leftarrow f_1, \bar{f}_2 \leftarrow f_2$  and  $\mathcal{D}_{\text{model}} = \emptyset$
- 2: **for**  $k = 0, \dots, K - 1$  **do**
- 3: Sample minibatch  $\mathcal{D}_{\text{real}}^{\text{mini}}$  from dataset  $\mathcal{D}_{\text{real}}$  and minibatch  $\mathcal{D}_{\text{model}}^{\text{mini}}$  from dataset  $\mathcal{D}_{\text{model}}$ .
- 4: Construct transition tuples using model predictions

$$\mathcal{D}_M := \{(s, a, r_M, s'_M) : r_M = R_M(s, a), s'_M \sim P_M(\cdot | s, a), (s, a) \in \mathcal{D}_{\text{real}}^{\text{mini}} \cup \mathcal{D}_{\text{model}}^{\text{mini}}\}$$

- 5: Update the adversary networks; for  $i = 1, 2$ ,

$$l^{\text{adversary}}(f, M) := \mathcal{L}_{\mathcal{D}_M}(f, \pi, \pi_{\text{ref}}) + \beta \left( \mathcal{E}_{\mathcal{D}_M}^w(f, M, \pi) + \lambda \mathcal{E}_{\mathcal{D}_{\text{model}}^{\text{mini}}}(M) \right) \quad (4)$$

$$M \leftarrow M - \eta_{\text{fast}} (\nabla_M l^{\text{adversary}}(f_1, M) + \nabla_M l^{\text{adversary}}(f_2, M))$$

$$f_i \leftarrow \text{Proj}_{\mathcal{F}}(f_i - \eta_{\text{fast}} \nabla_{f_i} l^{\text{adversary}}(f_i, M)) \quad \text{and} \quad \bar{f}_i \leftarrow (1 - \tau) \bar{f}_i + \tau f_i$$

- 6: Update actor network with respect to the first critic network and the reference policy

$$l^{\text{actor}}(\pi) := -\mathcal{L}_{\mathcal{D}_M}(f_1, \pi, \pi_{\text{ref}}) \quad (5)$$

$$\pi \leftarrow \text{Proj}_{\Pi}(\pi - \eta_{\text{slow}} \nabla_{\pi} l^{\text{actor}}(\pi))$$

- 7: If  $k\%H = 0$ , then reset model state:  $\bar{\mathcal{S}}_\pi \leftarrow \{s \in \mathcal{D}_{\text{real}}^{\text{mini}}\}$  and  $\bar{\mathcal{S}}_{\pi_{\text{ref}}} \leftarrow \{s \in \mathcal{D}_{\text{real}}^{\text{mini}}\}$
- 8: Query the MDP model to expand  $\mathcal{D}_{\text{model}}$  and update model state

$$\bar{A}_\pi := \{a : a \sim \pi(s), s \in \bar{\mathcal{S}}_\pi\} \quad \text{and} \quad \bar{A}_{\pi_{\text{ref}}} := \{a : a \sim \pi_{\text{ref}}(s), s \in \bar{\mathcal{S}}_{\pi_{\text{ref}}}\}$$

$$\mathcal{D}_{\text{model}} := \mathcal{D}_{\text{model}} \cup \{\bar{\mathcal{S}}_\pi, \bar{A}_\pi\} \cup \{\bar{\mathcal{S}}_{\pi_{\text{ref}}}, \bar{A}_{\pi_{\text{ref}}}\}$$

$$\bar{\mathcal{S}}_\pi \leftarrow \{s' : s' \sim \text{detach}(P_M(\cdot | s, a)), s \in \bar{\mathcal{S}}_\pi, a \in \bar{A}_\pi\}$$

$$\bar{\mathcal{S}}_{\pi_{\text{ref}}} \leftarrow \{s' : s' \sim \text{detach}(P_M(\cdot | s, a)), s \in \bar{\mathcal{S}}_{\pi_{\text{ref}}}, a \in \bar{A}_{\pi_{\text{ref}}}\}$$

- 9: **end for**

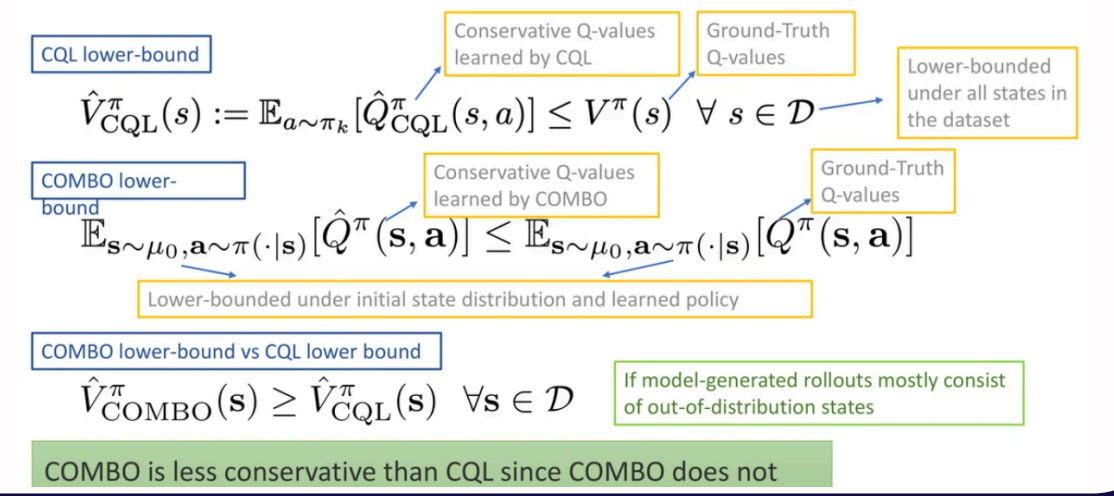
## COMBO

**Key Idea:** Augment the offline dataset with model rollouts while addressing distribution shift by learning lower-bounded Q-values without uncertainty quantification

$$\begin{aligned} \hat{Q}^{k+1} &\leftarrow \arg \min_Q \beta \left( \mathbb{E}_{\mathbf{s}, \mathbf{a} \sim \rho(\mathbf{s}, \mathbf{a})} [Q(\mathbf{s}, \mathbf{a})] - \mathbb{E}_{\mathbf{s}, \mathbf{a} \sim \mathcal{D}} [Q(\mathbf{s}, \mathbf{a})] \right) \\ &\quad + \frac{1}{2} \mathbb{E}_{\mathbf{s}, \mathbf{a}, \mathbf{s}' \sim d_f} \left[ \left( Q(\mathbf{s}, \mathbf{a}) - \hat{B}^\pi \hat{Q}^k(\mathbf{s}, \mathbf{a}) \right)^2 \right]. \end{aligned}$$

$$\pi' \leftarrow \arg \max_{\pi} \mathbb{E}_{\mathbf{s} \sim \rho, \mathbf{a} \sim \pi(\cdot | \mathbf{s})} \left[ \hat{Q}^\pi(\mathbf{s}, \mathbf{a}) \right].$$

### Theoretical Guarantees of COMBO



### Algorithm 1 COMBO: Conservative Model Based Offline Policy Optimization

**Require:** Offline dataset  $\mathcal{D}$ , rollout distribution  $\mu(\cdot | \mathbf{s})$ , learned dynamics model  $\hat{T}_\theta$ , initialized policy and critic  $\pi_\phi$  and  $Q_\psi$ .

- 1: Train the probabilistic dynamics model  $\hat{T}_\theta(s', r | s, \mathbf{a}) = \mathcal{N}(\mu_\theta(s, \mathbf{a}), \Sigma_\theta(s, \mathbf{a}))$  on  $\mathcal{D}$ .
- 2: Initialize the replay buffer  $\mathcal{D}_{\text{model}} \leftarrow \emptyset$ .
- 3: **for**  $i = 1, 2, 3, \dots$  **do**
- 4: Collect model rollouts by sampling from  $\mu$  and  $\hat{T}_\theta$  starting from states in  $\mathcal{D}$ . Add model rollouts to  $\mathcal{D}_{\text{model}}$ .
- 5: Conservatively evaluate  $\pi_\phi^i$  by repeatedly solving eq. 2 to obtain  $\hat{Q}_\psi^{\pi_\phi^i}$  using samples from  $\mathcal{D} \cup \mathcal{D}_{\text{model}}$ .
- 6: Improve policy under state marginal of  $d_f$  by solving eq. 3 to obtain  $\pi_\phi^{i+1}$ .
- 7: **end for**

## MOPO

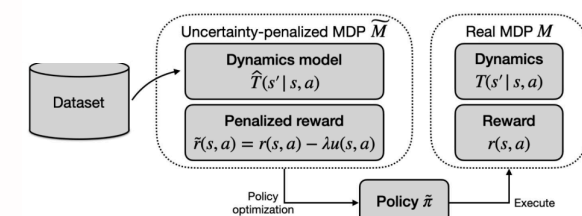
We optimize  $\pi$  in an *uncertainty-penalized* MDP  $\tilde{\mathcal{M}} = (S, \mathcal{A}, \hat{T}, \tilde{r}, \gamma)$ , where

$$\tilde{r}(s, a) = r(s, a) - \lambda u(s, a)$$

Key property: if  $u(s, a)$  is an *admissible error estimator*, i.e., upper bounds the error of  $\hat{T}(s' | s, a)$ , then for a particular choice of  $\lambda$ ,

$$\eta_{\mathcal{M}}(\pi) \geq \eta_{\tilde{\mathcal{M}}}(\pi)$$

where  $\eta_{\mathcal{M}}(\pi)$  and  $\eta_{\tilde{\mathcal{M}}}(\pi)$  denote the expected return of  $\pi$  in the MDPs  $\mathcal{M}$  and  $\tilde{\mathcal{M}}$ , respectively.



## ANALYSIS

### Standard Offline RL

(REF as **behavior cloned policy** on given dataset)

Algo	ARMOR	MoReL	MOPO	RAMBO	COMBO	ATAC	CQL	IQL
mujoco score	87.6	72.9	42.1	83.7	82.0	88.0	67.2	76.9

ARMOR is the only model-based method that doesn't rely on model ensembles!