

# A Unified Framework for Alternating Offline Model Training and Policy Learning

Shentao Yang<sup>1</sup>, Shujian Zhang<sup>1</sup>, Yihao Feng<sup>2</sup>, Mingyuan Zhou<sup>1</sup>

<sup>1</sup>The University of Texas at Austin, <sup>2</sup>Salesforce Research

## Proposed Method Sketch

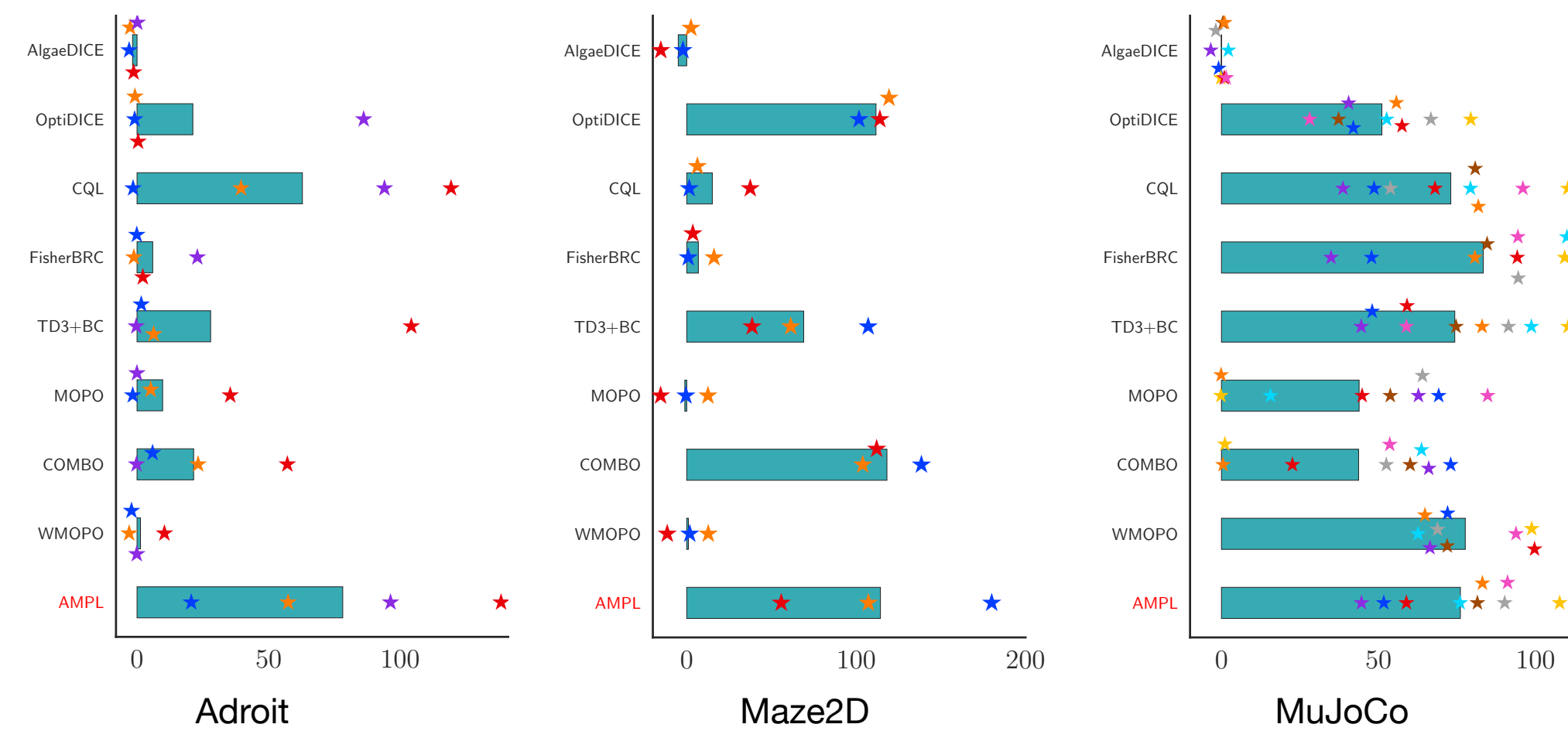
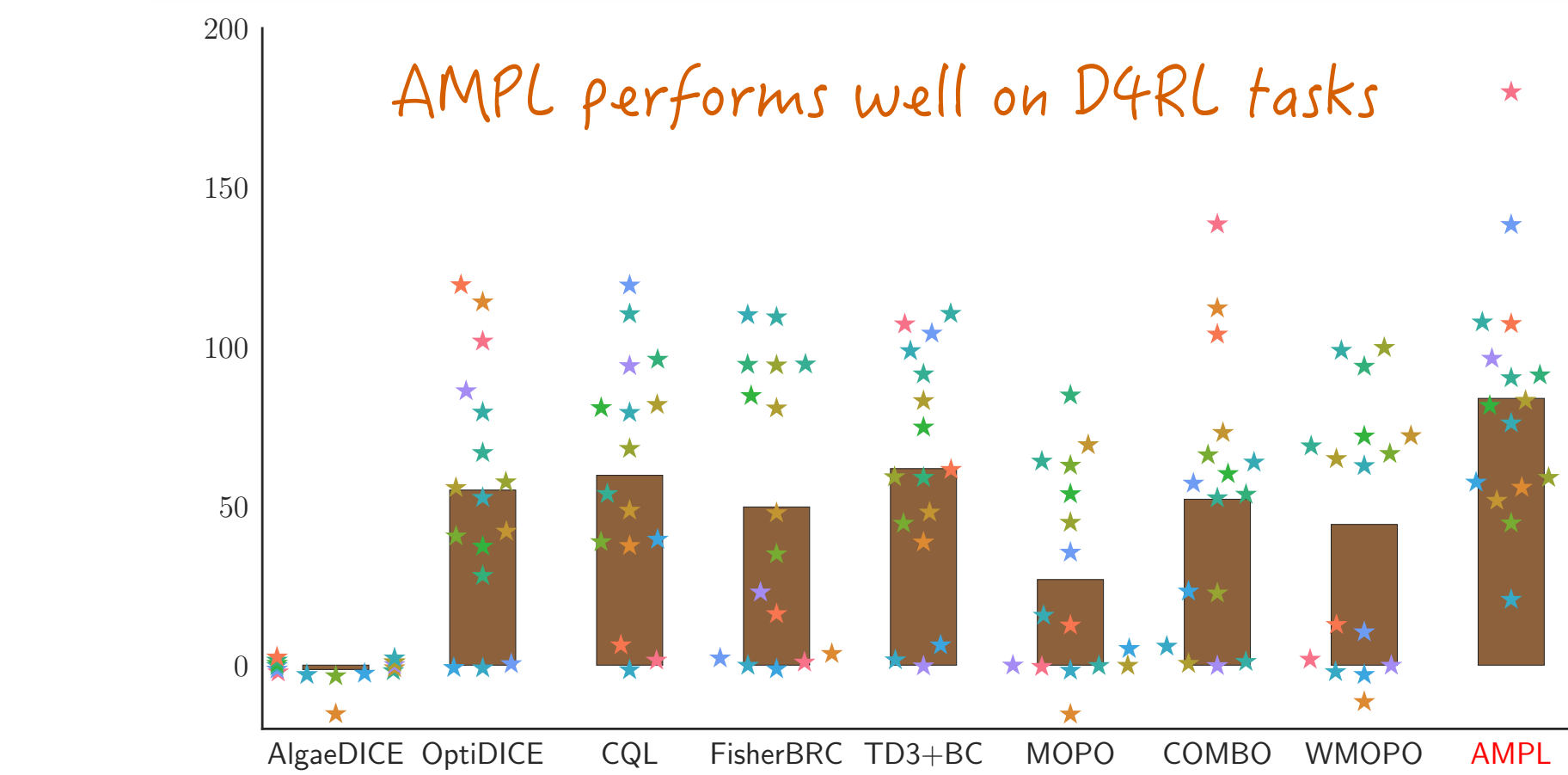
- **Motivation:** model training = MLE  $\neq$  improve policy = model usage.



$$\min_{\pi, \hat{P}} C \cdot \sqrt{D_{\pi}(P^*, \hat{P})} \geq |J(\pi, P^*) - J(\pi, \hat{P})|$$

- Jointly train and to minimize an upper bound of the evaluation error.
- Fixed ,  $\approx$  only on state-actions visited by .
- Fixed , optimize with a regularization based on .

## Results: AMPL



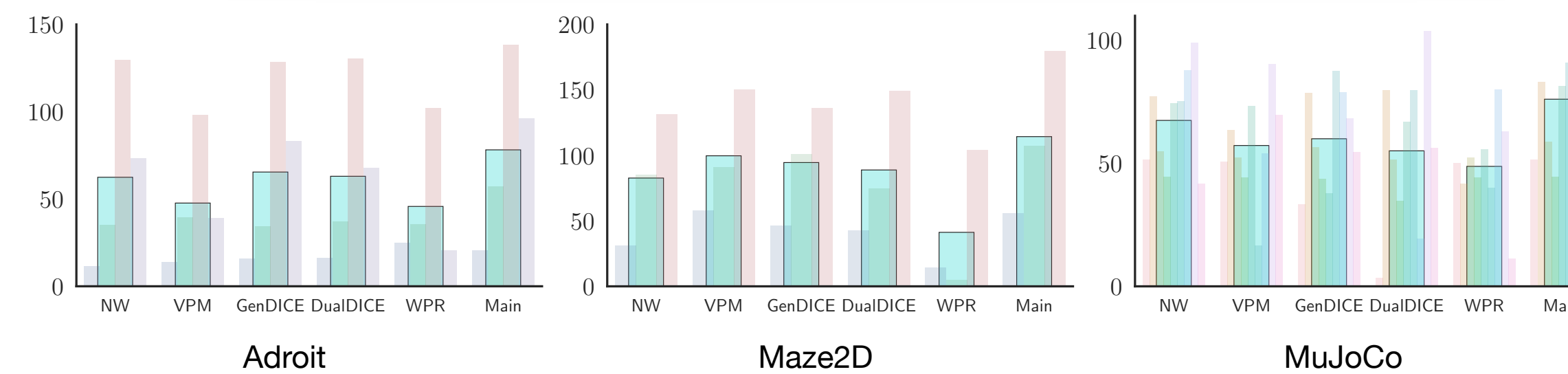
- Learn well on the MuJoCo datasets.
- Robust and good results on the challenging Adroit and Maze2D datasets.

## Background

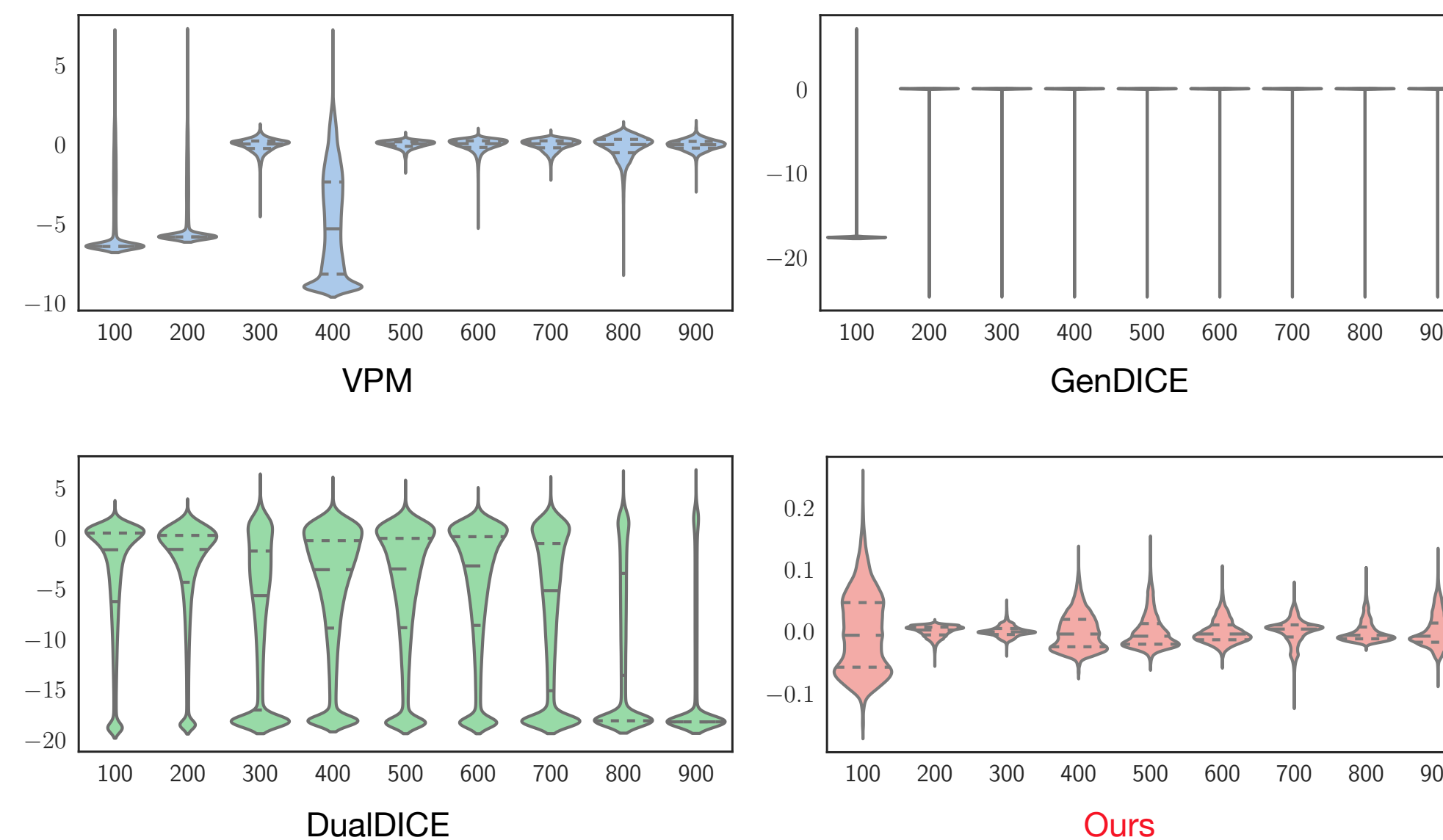


- Most offline MBRL: pre-train a fixed dynamic model on .
- Objective: MLE — “simply a mimic of the world.”
- Usage: improve the policy.
- **Objective mismatch:** model training  $\neq$  model usage.
- Especially when is limited and is hard to learn.

## Results: Ablation Study



- No Weights (NW): training only at the beginning using MLE.
- VPM, GenDICE, and DualDICE: other  $\omega(s, a)$  estimation methods.
- Can be unstable to provide good  $\omega(s, a)$  for training.



## Alternating Model-Policy Learning (AMPL)

- A tractable upper bound for the evaluation error

$$|J(\pi, P^*) - J(\pi, \hat{P})| \leq C \cdot \sqrt{D_{\pi}(P^*, \hat{P})}, \quad \text{with}$$

$$D_{\pi}(P^*, \hat{P}) \triangleq \mathbb{E}_{(s,a) \sim d_{\pi_b}^{P^*}} \left[ \omega(s, a) \text{KL} \left( P^*(s' | s, a) \pi_b(a' | s') \parallel \hat{P}(s' | s, a) \pi(a' | s') \right) \right],$$

- where  $\pi_b$  is the behavior policy ,  $d_{\pi_b, \gamma}^{P^*}$  is the offline-data distribution ,
- $\omega(s, a) \triangleq \frac{d_{\pi_b, \gamma}^{P^*}(s, a)}{d_{\pi_b, \gamma}(s, a)}$  is the density ratio between and visitation freq. of .
- Fix , train the model by  $\ell(\hat{P}) = -\mathbb{E}_{(s,a,s') \sim d_{\pi_b, \gamma}^{P^*}} \left[ \omega(s, a) \log \left\{ \hat{P}(s' | s, a) \right\} \right]$
- Given  $\omega(s, a)$ , a stable weighted MLE.

- Lower-bound of policy performance:  $J(\pi, \hat{P}) - C \cdot \sqrt{D_{\pi}(P^*, \hat{P})}$ .
- Fix , empirically helpful to construct the regularizer by:

- Removing  $\sqrt{\cdot}$ , applying a further relaxation before changing KL  $\rightarrow$  JSD
- $$D_{\pi}(P^*, \hat{P}) \leq C'' \cdot \text{KL} \left( P^*(s' | s, a) \pi_b(a' | s') d_{\pi_b, \gamma}^{P^*}(s, a) \parallel \hat{P}(s' | s, a) \pi(a' | s') d_{\pi_b, \gamma}^{P^*}(s) \pi(a | s) \right).$$

- **Density-ratio** training: Fixed-point-style simple MSE, ~~saddle-point optimization~~.

$$\mathbb{E}_{(s,a) \sim d_{\pi_b, \gamma}^{P^*}} \left[ \omega(s, a) \cdot Q_{\pi}^{\hat{P}}(s, a) \right] = \gamma \mathbb{E}_{(s,a,s') \sim d_{\pi_b, \gamma}^{P^*}} \left[ \omega(s, a) \cdot Q_{\pi}^{\hat{P}}(s', a') \right] + (1 - \gamma) \mathbb{E}_{s \sim \mu_0(\cdot)} \left[ Q_{\pi}^{\hat{P}}(s, a) \right].$$

$a' \sim \pi(\cdot | s')$        $a \sim \pi(\cdot | s)$

- Based on primal-dual relation between  $\omega(s, a)$  and Q-function in OPE.

## Summary

- **Goal:** close the mismatched model objectives in offline MBRL.
- **Method:** offline Alternating Model-Policy Learning.

QR code for the full paper!



QR code for the GitHub Repo!

