



# Deterministic Uncertainty Propagation for Improved Model-Based Offline Reinforcement Learning

**Abdullah Akgül, Manuel Haußmann, Melih Kandemir**

University of Southern Denmark

NeurIPS 2024

# The Scope

- Model-based Offline reinforcement learning
- Problems
  - ▶ Distributional shift
    - ★ Limited coverage on state-action space
  - ▶ Overestimation bias
    - ★ Errors due to policy search algorithms
    - ★ Yields suboptimal policies
  - ▶ Sampling and function approximation errors
    - ★ Further noise on training
    - ★ Decrease on learning speed
- PEssimistic Value Iteration (PEVI)<sup>1</sup>

---

<sup>1</sup>Jin et al., 2021. Is pessimism provably efficient for offline RL?

## Pessimistic Value Iteration

PEVI penalizes Bellman target estimation with the uncertainty on the predicted next state to minimize the suboptimality of a policy  $\pi$ :

$$\text{SubOpt}(\pi; s) \triangleq Q_{\pi^*}(s, \pi^*(s)) - Q_{\pi}(s, \pi(s))$$

for an initial state  $s$ .

## Pessimistic Value Iteration

PEVI penalizes Bellman target estimation with the uncertainty on the predicted next state to minimize the suboptimality of a policy  $\pi$ :

$$\text{SubOpt}(\pi; s) \triangleq Q_{\pi^*}(s, \pi^*(s)) - Q_{\pi}(s, \pi(s))$$

for an initial state  $s$ .

### Theorem (*Suboptimality of PEVI*)

For any  $\pi$  derived with PEVI that satisfies

$$\underbrace{\left| \mathbb{B}_{\pi} Q(s, a) - \widehat{\mathbb{B}}_{\pi} Q(s, a, s') \right|}_{\text{Bellman approximation error}} \leq \underbrace{\Gamma_{\widehat{\mathbb{P}}}^Q(s, a)}_{\text{Uncertainty quantifier}}, \quad \forall (s, a) \in \mathcal{S} \times \mathcal{A}$$

with probability at least  $1 - \delta$  for some error tolerance  $\delta \in (0, 1)$ , the following inequality holds:

$$\text{SubOpt}(\pi; s) \leq f(\Gamma_{\widehat{\mathbb{P}}}^Q, s, \pi^*).$$

# PEVI Approaches

- MOPO<sup>2</sup> penalizes via uncertainty on the next state
- MOBILE<sup>3</sup> penalizes via uncertainty on the Bellman target

Both approximate the Bellman target by evaluating with a sample  $s'$ .

---

<sup>2</sup>Yu et al., 2020. MOPO: Model-based offline policy optimization

<sup>3</sup>Sun et al., 2023. Model-Bellman inconsistency for model-based offline reinforcement learning

# PEVI Approaches

- MOPO<sup>2</sup> penalizes via uncertainty on the next state
- MOBILE<sup>3</sup> penalizes via uncertainty on the Bellman target

Both approximate the Bellman target by evaluating with a sample  $s'$ .

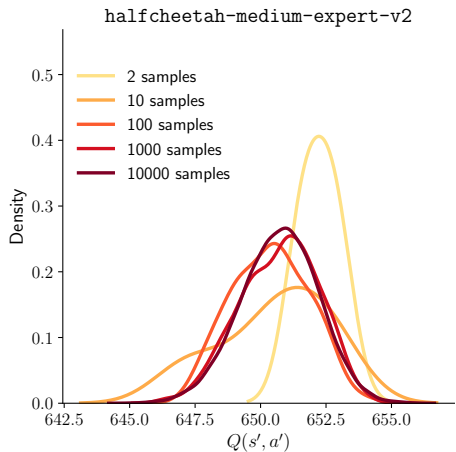
**1. Contribution:** We provide a suboptimality guarantee for sampling-based PEVI approaches.

---

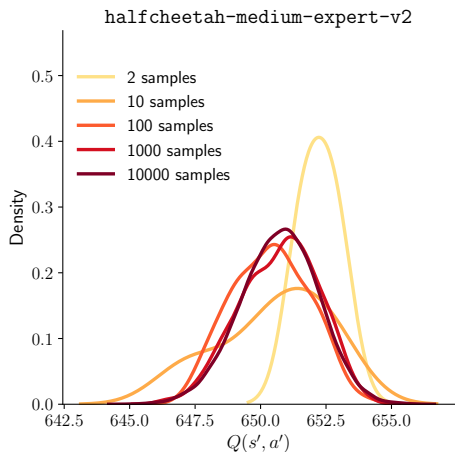
<sup>2</sup>Yu et al., 2020. MOPO: Model-based offline policy optimization

<sup>3</sup>Sun et al., 2023. Model-Bellman inconsistency for model-based offline reinforcement learning

# The Challenge: High Variance



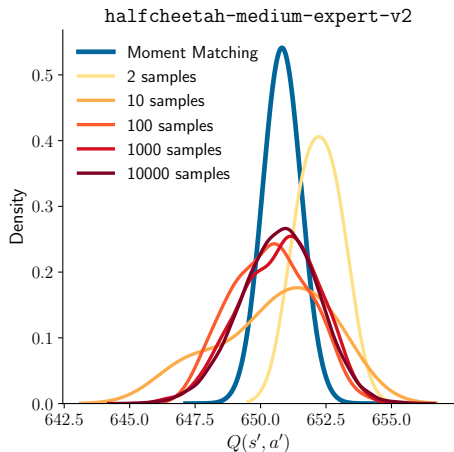
# The Challenge: High Variance



- Distorted gradient signals, delayed convergence
- Poor approximation of the first two moments of Bellman target
- Requirement of larger confidence radii

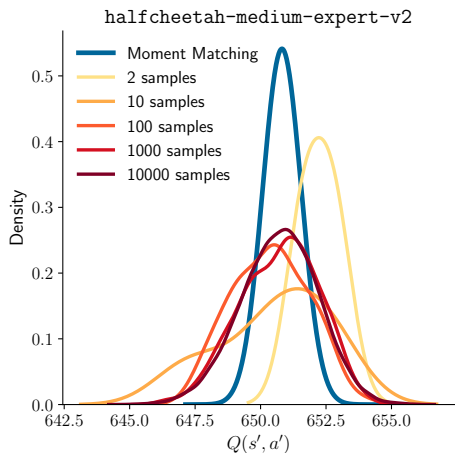


# The Challenge: High Variance



- Distorted gradient signals, delayed convergence
- Poor approximation of the first two moments of Bellman target
- Requirement of larger confidence radii

# The Challenge: High Variance



- Distorted gradient signals, delayed convergence
- Poor approximation of the first two moments of Bellman target
- Requirement of larger confidence radii

## 2. Contribution: Moment matching.

# The Solution

## MOMBO: Moment Matching Offline Model-Based Policy Optimization

- Deterministic uncertainty propagation
  - ▶ Propagating first two moments of uncertain input through a value function
    - ★ Neural network
- Lower confidence bound on the estimation of Bellman target

# The Solution

## MOMBO: Moment Matching Offline Model-Based Policy Optimization

- Deterministic uncertainty propagation
  - ▶ Propagating first two moments of uncertain input through a value function
    - ★ Neural network
- Lower confidence bound on the estimation of Bellman target
  
- **3. Contribution:** Suboptimality bound for moment matching
  - ▶ Tighter bound
  - ▶ Provably more efficient

# The Experiments

## Performance Evaluation

Dataset Type	Environment	NORMALIZED REWARD ( $\uparrow$ )			AULC ( $\uparrow$ )		
		MOPO	MOBILE	MOMBO	MOPO	MOBILE	MOMBO
random	halfcheetah	37.2 $\pm$ 1.6	41.2 $\pm$ 1.1	<b>43.6<math>\pm</math>1.1</b>	36.3 $\pm$ 1.0	39.5 $\pm$ 1.2	<b>41.4<math>\pm</math>1.0</b>
	hopper	<b>31.7<math>\pm</math>0.1</b>	31.3 $\pm$ 0.1	25.4 $\pm$ 10.2 $^\dagger$	<b>28.6<math>\pm</math>1.4</b>	23.6 $\pm$ 3.7	17.3 $\pm$ 1.3
	walker2d	8.2 $\pm$ 5.6	<b>22.1<math>\pm</math>0.9</b>	<u>21.5<math>\pm</math>0.1</u>	5.4 $\pm$ 3.2	18.0 $\pm$ 0.4	<b>19.2<math>\pm</math>0.5</b>
	Average	25.7	<b>31.5</b>	30.2	23.4	<b>27.1</b>	26.0
medium	halfcheetah	72.4 $\pm$ 4.2	<u>75.8<math>\pm</math>0.8</u>	<b>76.1<math>\pm</math>0.8</b>	70.9 $\pm$ 2.0	<u>72.1<math>\pm</math>1.0</u>	<b>73.0<math>\pm</math>0.9</b>
	hopper	62.8 $\pm$ 38.1	103.6 $\pm$ 1.0	<b>104.2<math>\pm</math>0.5</b>	37.0 $\pm$ 15.3	82.2 $\pm$ 7.3	<b>95.9<math>\pm</math>2.5</b>
	walker2d	85.4 $\pm$ 2.9	<b>88.3<math>\pm</math>2.5</b>	<u>86.4<math>\pm</math>1.2</u>	77.6 $\pm$ 1.3	79.0 $\pm$ 1.3	<b>84.0<math>\pm</math>1.1</b>
	Average	73.6	<b>89.3</b>	88.9	61.8	77.8	<b>84.3</b>
medium-replay	halfcheetah	<b>72.1<math>\pm</math>3.8</b>	<u>71.9<math>\pm</math>3.2</u>	<u>72.0<math>\pm</math>4.3</u>	<u>68.4<math>\pm</math>4.7</u>	<u>67.9<math>\pm</math>2.8</u>	<b>68.7<math>\pm</math>3.9</b>
	hopper	92.7 $\pm$ 20.7	<b>105.1<math>\pm</math>1.3</b>	<u>104.8<math>\pm</math>1.0</u>	81.7 $\pm$ 4.6	78.7 $\pm$ 4.0	<b>87.3<math>\pm</math>2.0</b>
	walker2d	85.9 $\pm$ 5.3	<b>90.5<math>\pm</math>1.7</b>	<u>89.6<math>\pm</math>3.8</u>	65.3 $\pm$ 12.7	<u>79.9<math>\pm</math>4.3</u>	<b>80.8<math>\pm</math>5.6</b>
	Average	83.4	<b>89.2</b>	88.8	72.4	75.5	<b>78.9</b>
medium-expert	halfcheetah	83.6 $\pm$ 12.5	100.9 $\pm$ 1.5	<b>103.3<math>\pm</math>0.8</b>	77.1 $\pm$ 4.0	<u>94.5<math>\pm</math>1.8</u>	<b>95.2<math>\pm</math>0.7</b>
	hopper	74.9 $\pm$ 44.2	<u>112.5<math>\pm</math>0.2</u>	<b>112.6<math>\pm</math>0.3</b>	55.6 $\pm$ 17.3	<u>82.7<math>\pm</math>7.3</u>	<b>84.3<math>\pm</math>4.7</b>
	walker2d	108.2 $\pm$ 4.3	<b>114.5<math>\pm</math>2.2</b>	<u>113.9<math>\pm</math>0.9</u>	88.3 $\pm$ 6.3	94.3 $\pm$ 0.9	<b>98.9<math>\pm</math>3.3</b>
	Average	88.9	109.3	<b>109.9</b>	73.6	90.5	<b>92.8</b>
Average Score		67.6	<b>79.8</b>	79.5	57.5	67.7	<b>70.5</b>
Average Ranking		2.7	<b>1.7</b>	<b>1.7</b>	2.7	2.2	<b>1.2</b>

$^\dagger$  High standard deviation due to failure in one repetition, which can be mitigated by increasing  $\beta$ . Median result: 31.3

# Conclusion

We introduce MOMBO

- has faster convergence and more stable training
- provides a competitive final performance
- estimates Bellman target more precisely

# Conclusion

We introduce MOMBO

- has faster convergence and more stable training
- provides a competitive final performance
- estimates Bellman target more precisely

