

# Lecture 7: Offline RL

# CS885 Reinforcement Learning

2025-01-28

Complementary readings:

Levine, Kumar, Tucker, Fu (2021) Offline reinforcement learning: Tutorial, review, and perspectives on open problems, *arxiv*.

Kumar, Zhou, Tucker, Levine (2020) Conservative Q-Learning for Offline Reinforcement Learning, *NeurIPS*.

Pascal Poupart

David R. Cheriton School of Computer Science

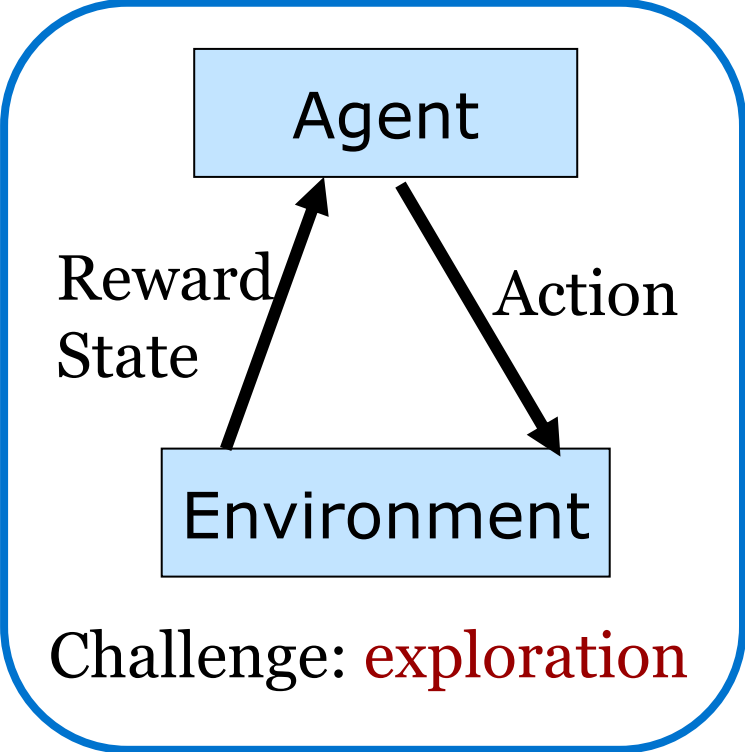


# Outline

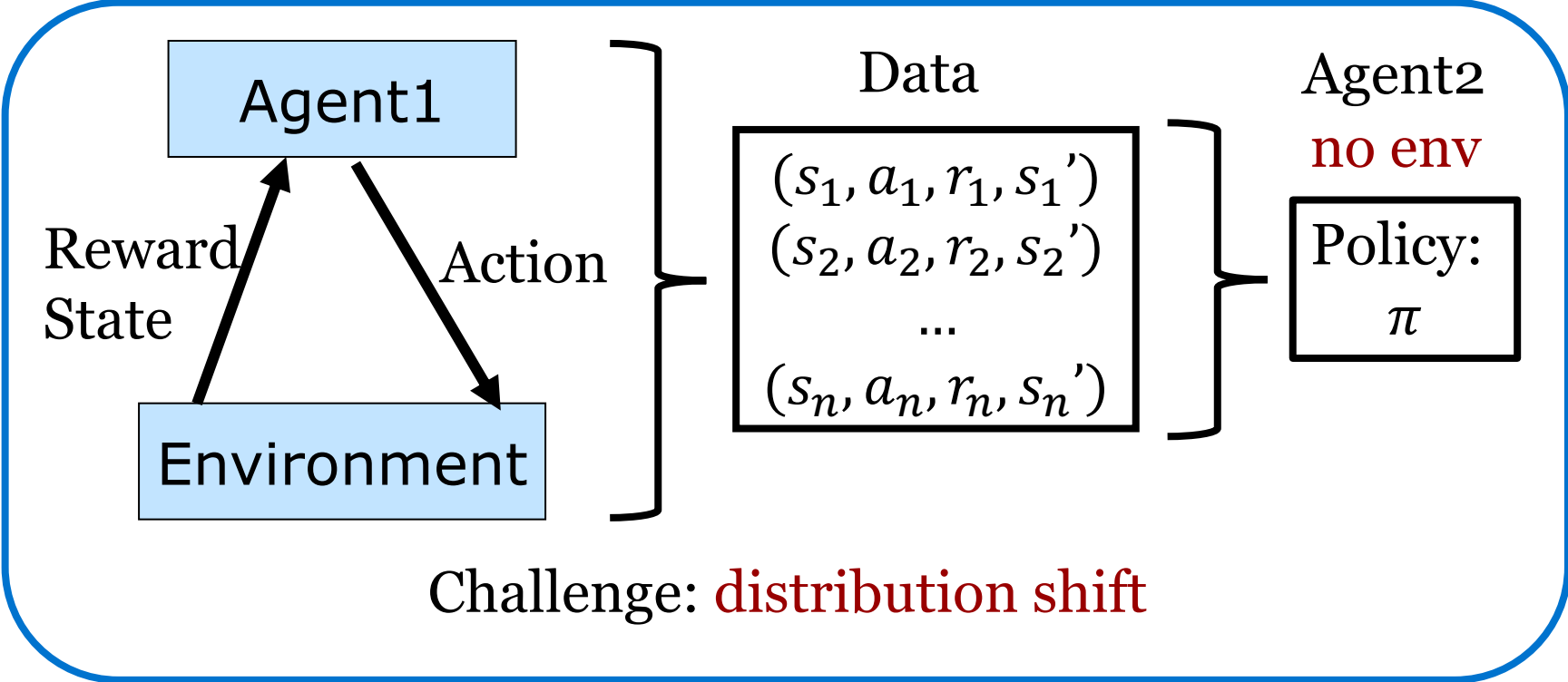
- Can we optimize a policy without interacting with the environment (i.e., learn from previously saved data)?
- Offline RL (also known as batch RL)
  - Conservative Q-Learning
  - Conservative Soft Q-learning
  - Conservative Soft Actor Critic (SAC)

# Reinforcement Learning

## Online RL

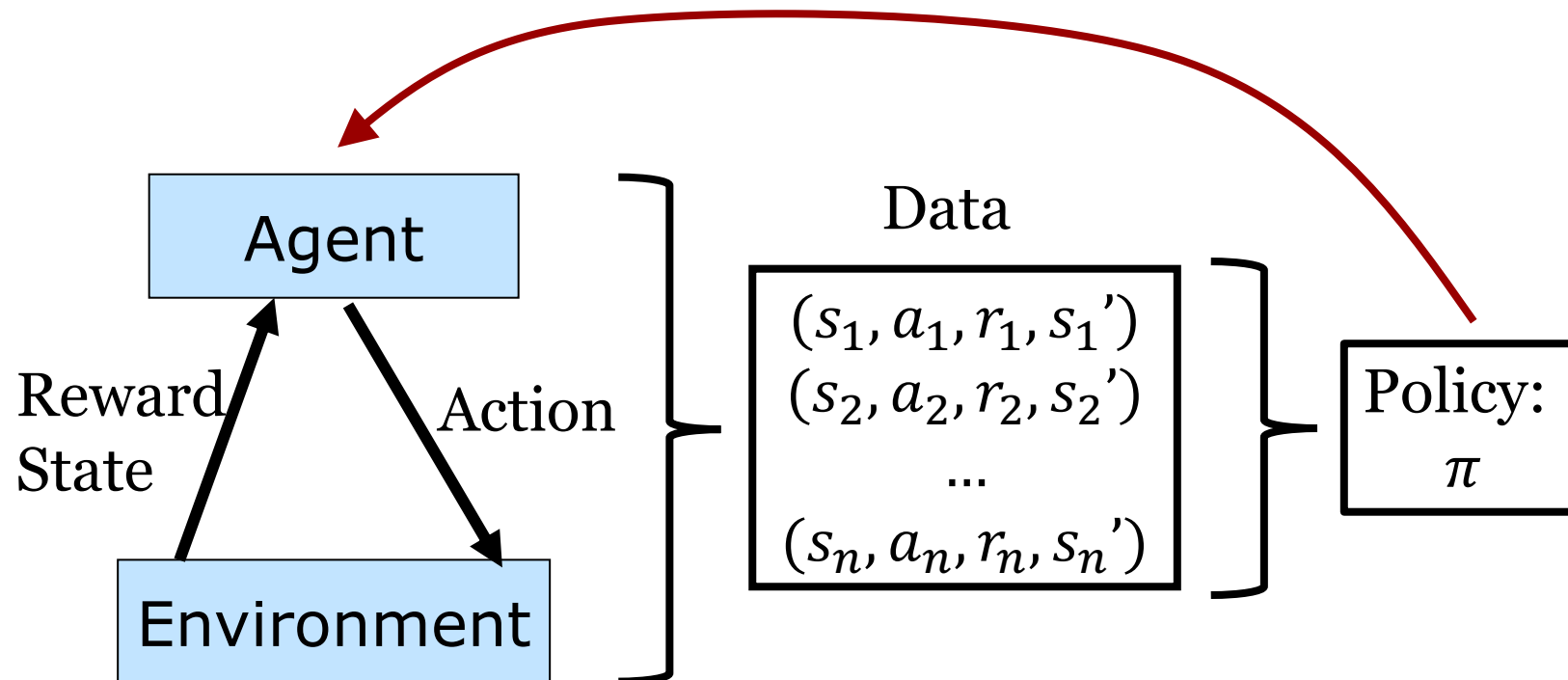


## Offline RL



# Off-Policy RL

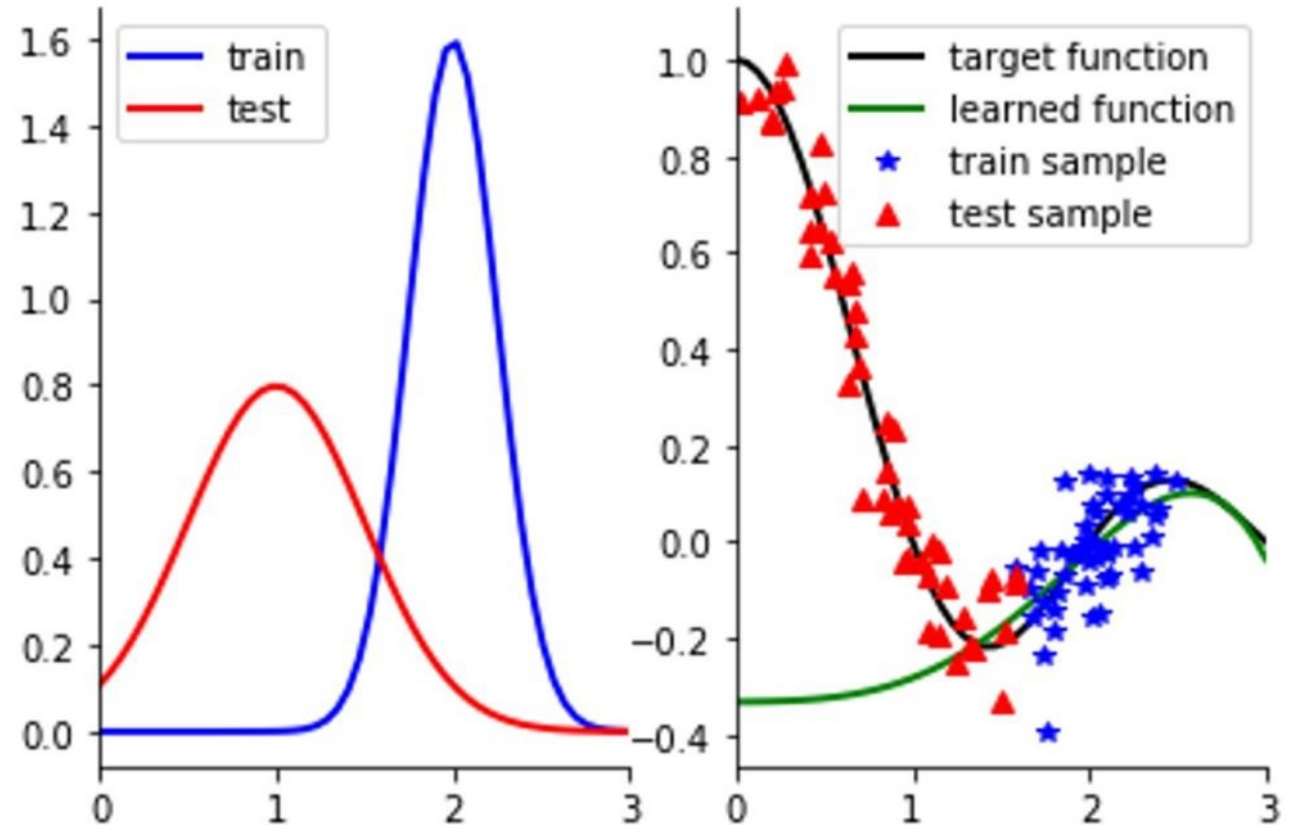
- Form of online RL since agent can experiment with its policy in environment



# Distribution Shift

- Train distribution different from test distribution
- In RL: data generated by  $\pi$ , but goal is to learn improved  $\pi'$
- Challenge: may choose actions with overestimated Q-values

From Christyn Zehnder (in-Q-Tel)



# Offline RL Techniques

- Importance Sampling
- Policy constraints
- Penalty methods
  - Conservative Q-Learning, conservative Soft Actor Critic
- Model-based RL

# Off-Policy Evaluation by Q-learning

- Let  $\pi_\beta(a|s)$  be a behaviour policy to collect  $D = \{(s, a, r, s')\}$
- We can evaluate a different policy  $\pi$  by off-policy Q-learning:

$$Q^\pi = \operatorname{argmin}_Q E_{(s,a,r,s') \sim D} \left[ \left( r + \gamma E_{a' \sim \pi(a'|s')} [Q(s', a')] - Q(s, a) \right)^2 \right]$$

- Some Q-values underestimated and others overestimated
- Greedy policy improvement:  $\pi_{k+1}(s) \leftarrow \operatorname{argmax}_a Q^{\pi_k}(s, a) \forall s$ 
  - **Problem: select actions with overestimated Q-values**

# Conservative Off-Policy Evaluation

- Introduce a penalty term

$$\hat{Q}^\pi = \operatorname{argmin}_Q \eta E_{s \sim D, a \sim \pi(a|s)}[Q(s, a)] + E_{(s, a, r, s') \sim D} \left[ \left( r + \gamma E_{a' \sim \pi(a'|s')} Q(s', a') - Q(s, a) \right)^2 \right]$$

where  $\eta$ : weight that determines importance of penalty

- Let  $\operatorname{support}(\pi) = \{(s, a) \mid \pi \text{ reaches } (s, a) \text{ with non-zero probability}\}$

**Theorem:** If  $\operatorname{support}(\pi) \subseteq \operatorname{support}(\pi_\beta)$ , then for sufficiently large  $\eta$ ,

(from Kumar et al. 2020)  $\hat{Q}^\pi(s, a) \leq Q^\pi(s, a) \quad \forall s \in D, a$



# Improved Bound

- Remove  $E_{s,a \sim D}[Q(s, a)]$  from penalty term

$$\tilde{Q}^\pi = \operatorname{argmin}_Q \eta \left( E_{s \sim D, a \sim \pi(a|s)}[Q(s, a)] - E_{s,a \sim D}[Q(s, a)] \right) \\ + E_{(s,a,r,s') \sim D} \left[ \left( r + \gamma E_{a' \sim \pi(a'|s')} Q(s', a') - Q(s, a) \right)^2 \right]$$

- We cannot guarantee that  $\tilde{Q}^\pi(s, a) \leq Q^\pi(s, a) \forall s \in D, a$  for sufficiently large  $\eta$
- Let  $V^\pi(s) = E_{a \sim \pi(a|s)} Q^\pi(s, a)$

**Theorem:** If  $\operatorname{support}(\pi) \subseteq \operatorname{support}(\pi_\beta)$ , then for sufficiently large  $\eta$ ,

(from Kumar et al. 2020) 
$$\tilde{V}^\pi(s) \leq V^\pi(s) \quad \forall s \in D$$

# Conservative Q-learning

- Idea: let  $\pi$  be the greedy policy:  $\pi(s) = \operatorname{argmax}_a Q(s, a)$

$$\begin{aligned} \tilde{Q}^* = \operatorname{argmin}_Q \eta & \left( E_{s \sim D} \left[ \max_a Q(s, a) \right] - E_{s, a \sim D} [Q(s, a)] \right) \\ & + E_{(s, a, r, s') \sim D} \left[ \left( r + \gamma \max_{a'} Q(s', a') - Q(s, a) \right)^2 \right] \end{aligned}$$

# Conservative Q-Learning

Load fixed buffer of experiences

Initialize weights  $\mathbf{w}$  and  $\bar{\mathbf{w}}$  at random in  $[-1,1]$

Loop

Sample minibatch of  $n$  experiences from buffer

$$\text{Bellman error: } Err(\mathbf{w}) = \frac{1}{n} \sum_{(s,a,r,s') \in \text{minibatch}} \left[ \left( Q_{\mathbf{w}}(s, a) - r - \gamma \max_{a'} Q_{\bar{\mathbf{w}}}(s', a') \right)^2 \right]$$

$$\text{Penalty: } Penalty(\mathbf{w}) = \frac{1}{n} \sum_{(s,a) \in \text{minibatch}} \left[ \max_{\hat{a}} Q_{\mathbf{w}}(s, \hat{a}) - Q_{\mathbf{w}}(s, a) \right]$$

$$\text{Update weights: } \mathbf{w} \leftarrow \mathbf{w} - \alpha \left( \frac{\partial Err}{\partial \mathbf{w}} + \eta \frac{\partial Penalty}{\partial \mathbf{w}} \right)$$

Every  $c$  steps, update target:  $\bar{\mathbf{w}} \leftarrow \mathbf{w}$

# Conservative Soft Q-Learning

Load fixed buffer of experiences

Initialize weights  $\mathbf{w}$  and  $\bar{\mathbf{w}}$  at random in  $[-1,1]$

Loop

Sample minibatch of  $n$  experiences from buffer

$$\text{Bellman error: } Err(\mathbf{w}) = \frac{1}{n} \sum_{(s,a,r,s') \in \text{minibatch}} \left[ \left( Q_{\mathbf{w}}(s, a) - r - \gamma \tilde{\max}_{a'} Q_{\bar{\mathbf{w}}}(s', a') \right)^2 \right]$$

$$\text{Penalty: } Penalty(\mathbf{w}) = \frac{1}{n} \sum_{(s,a) \in \text{minibatch}} \left[ \tilde{\max}_{\hat{a}} Q_{\mathbf{w}}(s, \hat{a}) - Q_{\mathbf{w}}(s, a) \right]$$

$$\text{Update weights: } \mathbf{w} \leftarrow \mathbf{w} - \alpha \left( \frac{\partial Err}{\partial \mathbf{w}} + \eta \frac{\partial Penalty}{\partial \mathbf{w}} \right)$$

Every  $c$  steps, update target:  $\bar{\mathbf{w}} \leftarrow \mathbf{w}$

# Conservative Soft Actor Critic (SAC)

Load fixed buffer of experiences

Initialize weights  $\mathbf{w}$ ,  $\bar{\mathbf{w}}$  and  $\boldsymbol{\theta}$  at random in  $[-1,1]$

Loop

Sample minibatch of  $n$  experiences from buffer

For each experience  $(s, a, r, s')$  in minibatch, sample  $a' \sim \pi_{\boldsymbol{\theta}}(a'|s')$

Bellman error:  $Err(\mathbf{w}) = \frac{1}{n} \sum_{(s,a,r,s',a') \in minibatch} \left[ (Q_{\mathbf{w}}(s, a) - r - \gamma [Q_{\bar{\mathbf{w}}}(s', a') + \lambda H(\pi_{\boldsymbol{\theta}}(\cdot | s'))])^2 \right]$

Penalty:  $Penalty(\mathbf{w}) = \frac{1}{n} \sum_{(s,a) \in minibatch} [\widetilde{\max}_{\hat{a}} Q_{\mathbf{w}}(s, \hat{a}) - Q_{\mathbf{w}}(s, a)]$

Q-function update:  $\mathbf{w} \leftarrow \mathbf{w} - \alpha \left( \frac{\partial Err}{\partial \mathbf{w}} + \eta \frac{\partial Penalty}{\partial \mathbf{w}} \right)$

Policy update: Update policy:  $\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} - \alpha \frac{\partial KL(\pi_{\boldsymbol{\theta}} | softmax(Q_{\bar{\mathbf{w}}}/\lambda))}{\partial \boldsymbol{\theta}}$

Every  $c$  steps, update target:  $\bar{\mathbf{w}} \leftarrow \mathbf{w}$

# Empirical Evaluation

Conservative  
SAC

- Kumar et al.  
(NeurIPS-2020)

Task Name	SAC	BC	BEAR	BRAC-p	BRAC-v	CQL( $\mathcal{H}$ )
halfcheetah-random	30.5	2.1	25.5	23.5	28.1	<b>35.4</b>
hopper-random	<b>11.3</b>	9.8	9.5	<b>11.1</b>	<b>12.0</b>	<b>10.8</b>
walker2d-random	4.1	1.6	<b>6.7</b>	0.8	0.5	<b>7.0</b>
halfcheetah-medium	-4.3	36.1	38.6	<b>44.0</b>	<b>45.5</b>	<b>44.4</b>
walker2d-medium	0.9	6.6	33.2	72.7	<b>81.3</b>	79.2
hopper-medium	0.8	29.0	47.6	31.2	32.3	<b>58.0</b>
halfcheetah-expert	-1.9	<b>107.0</b>	<b>108.2</b>	3.8	-1.1	104.8
hopper-expert	0.7	<b>109.0</b>	<b>110.3</b>	6.6	3.7	<b>109.9</b>
walker2d-expert	-0.3	125.7	106.1	-0.2	-0.0	<b>153.9</b>
halfcheetah-medium-expert	1.8	35.8	51.7	43.8	45.3	<b>62.4</b>
walker2d-medium-expert	1.9	11.3	10.8	-0.3	0.9	<b>98.7</b>
hopper-medium-expert	1.6	<b>111.9</b>	4.0	1.1	0.8	<b>111.0</b>
halfcheetah-random-expert	53.0	1.3	24.6	30.2	2.2	<b>92.5</b>
walker2d-random-expert	0.8	0.7	1.9	0.2	2.7	<b>91.1</b>
hopper-random-expert	5.6	10.1	10.1	5.8	11.1	<b>110.5</b>
halfcheetah-mixed	-2.4	38.4	36.2	<b>45.6</b>	<b>45.9</b>	<b>46.2</b>
hopper-mixed	3.5	11.8	25.3	0.7	0.8	<b>48.6</b>
walker2d-mixed	1.9	11.3	10.8	-0.3	0.9	<b>26.7</b>

Table 1: Performance of CQL( $\mathcal{H}$ ) and prior methods on gym domains from D4RL, on the normalized return metric, averaged over 4 seeds. Note that CQL performs similarly or better than the best prior method with simple datasets, and greatly outperforms prior methods with complex distributions (“-mixed”, “-random-expert”, “-medium-expert”).