



# Offline Reinforcement Learning

Nuowen Kan

Shanghai Jiao Tong University

May 2024





# Contents at a glance

1. Preliminary
2. Challenges of Offline RL
3. Solutions
  - a) Policy constraint methods
  - b) Value-regularized method
  - c) Implicit constraints/regularization
  - d) Model-based method
  - e) Conditional imitation method



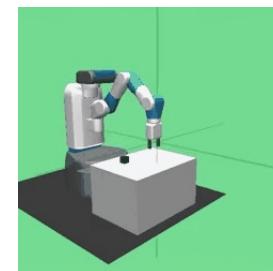
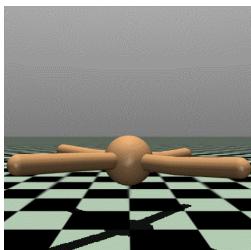
# Why offline RL?

- **Success of modern machine learning**  
Data + huge network



Can we develop data-driven RL methods?

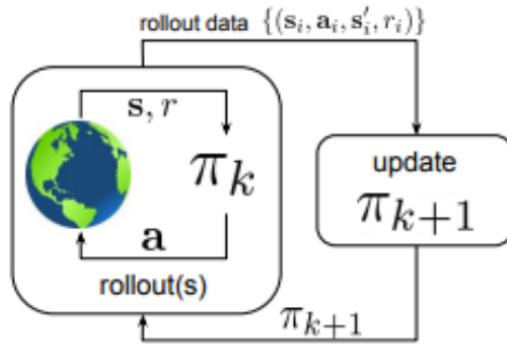
- **Training of RL**



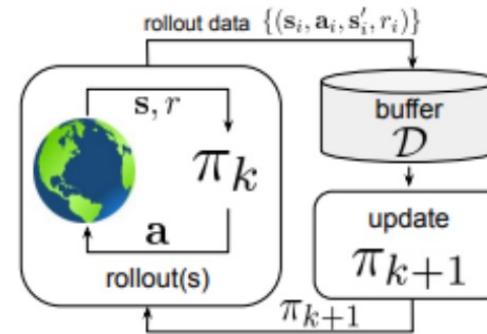


# On/off-Policy? On/off-Line?

## On-Policy



## Off-Policy



Reinforcement Learning with Online Interactions

## Online



Offline Reinforcement Learning

## Offline





- **Why, or when, might offline RL be more useful?**

leverage datasets collected by people, existing systems

online policy collection may be risky, unsafe

reuse previously collected data rather than recollecting

(e.g. previous experiments, projects, robots, institutions)

Note: A blend of **offline** then **online RL** is also possible!



- **Offline RL terminology**

$$\mathcal{D} = \{(s_i, a_i, s'_i, r_i)\}$$

Offline dataset sampled from some unknown policy  $\pi_\beta(a|s)$

$$\pi_\beta(a|s)$$

Behavior/sampling policy

$$a \sim \pi_\beta(a|s)$$

$$s \sim d^{\pi_\beta}(s)$$

$$s' \sim p(s'|s, a)$$

$$r \leftarrow r(s, a)$$



Sampled trajectories

Objective

$$\max_{\pi} \sum_{t=0}^T E_{s_t \sim d^\pi(s), a_t \sim \pi(a|s)} [\gamma^t r(s_t, a_t)]$$



- **Where does the data come from?**

human collected data

data from a hand-designed system / controller

data from previous RL run(s)

a mixture of sources



# Contents at a glance

1. Preliminary
2. Challenges of Offline RL
3. Solutions
  - a) Policy constraint methods
  - b) Value-regularized method
  - c) Implicit constraints/regularization
  - d) Model-based method
  - e) Conditional imitation method



# Challenges in offline RL

- **Exploration is impossible**
- **Counterfactual query during offline training**
  - Must act different to learn a better policy
  - Against I.I.D assumption
- **Distributional shift**
  - Distribution mismatch: In offline training, RL policy is trained under distribution  $\mathcal{D}$ , while it's evaluated on another.
  - Causes: new policy is trained by maximizing the expected return, which changes in visited states during testing.

<https://research.google/blog/tackling-open-challenges-in-offline-reinforcement-learning/>



- Can we just use off-policy algorithms?

Q-learning objective

$$\sum_{(\mathbf{s}, \mathbf{a}, \mathbf{s}') \sim \mathcal{D}} \| Q(\mathbf{s}, \mathbf{a}) - (r(\mathbf{s}, \mathbf{a}) + \gamma \max_{\mathbf{a}'} Q(\mathbf{s}', \mathbf{a}')) \|^2$$

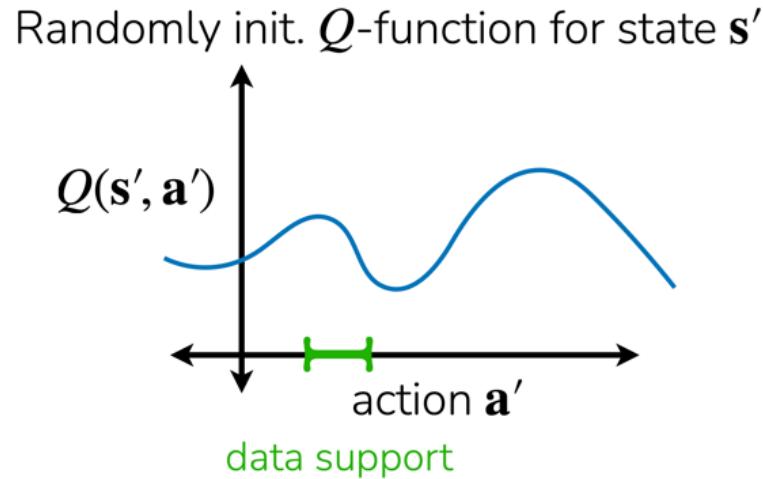
What happens if you optimize this **using a static dataset?**



Evaluating  $Q$  on actions  $\mathbf{a}'$  **not in the dataset**



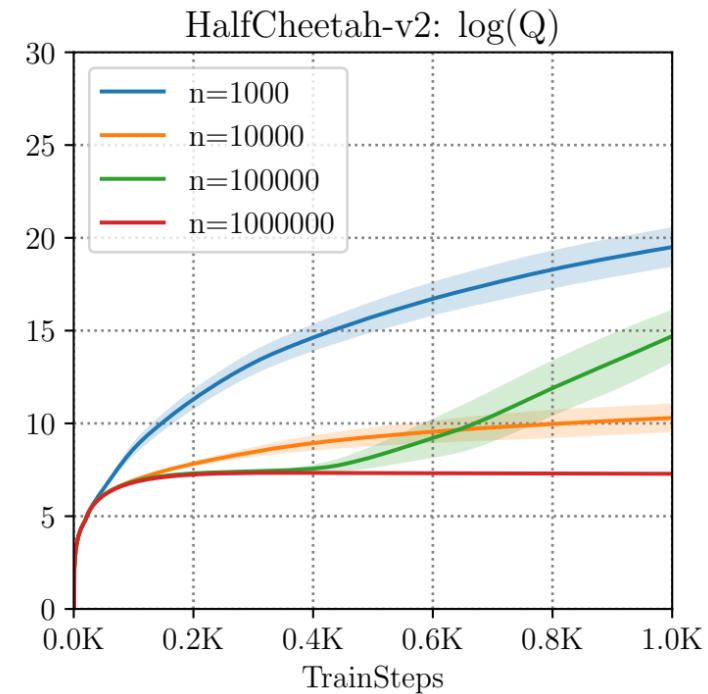
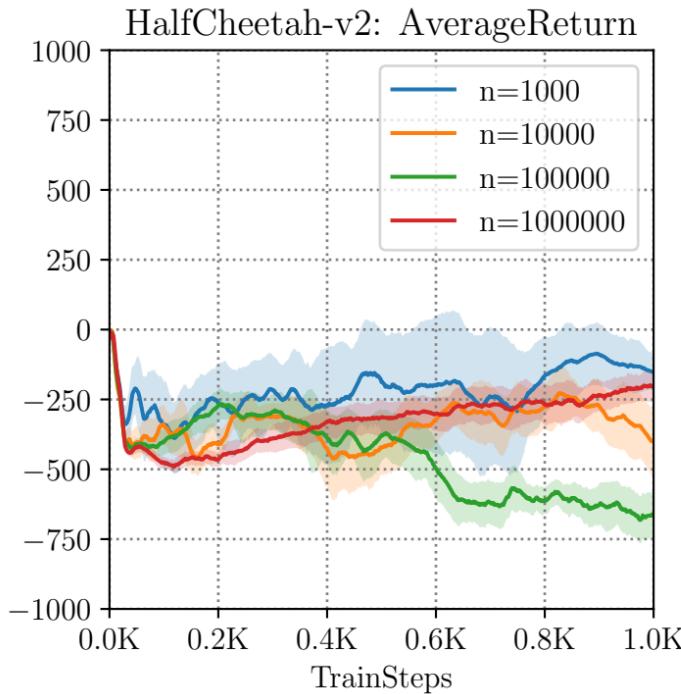
- Suffer from distribution shift



- Q-function will be **unreliable** on Out-of-distributional (OOD) actions
- **max** will seek out actions where Q-function is **over-optimistic**
- After values propagate, Q-values will become substantially **over-estimated**



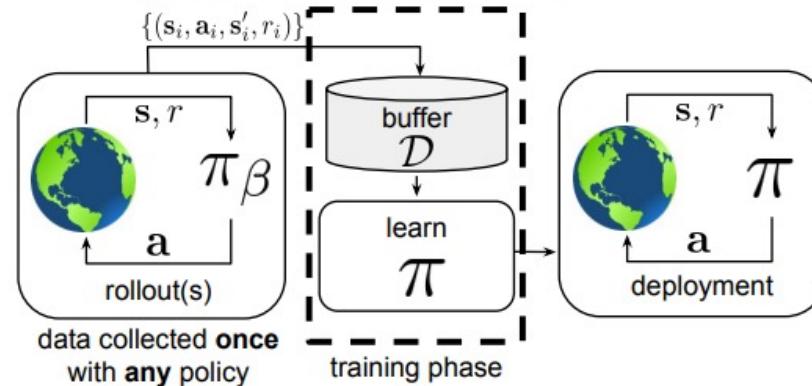
- Suffer from distribution shift



How to mitigate overestimation in offline RL



- Suffer from distribution shift



- Estimate the expectation of a different distribution

## Important Sampling

Unbiased estimated but  
large variation

$$\begin{aligned}\mathbb{E}_{X \sim P}[f(X)] &= \sum P(X)f(X) \\ &= \sum Q(X)\frac{P(X)}{Q(X)}f(X) \\ &= \mathbb{E}_{X \sim Q}\left[\frac{P(X)}{Q(X)}f(X)\right]\end{aligned}$$



# Contents at a glance

1. Preliminary
2. Challenges of Offline RL
3. Solutions
  - a) Policy constraint methods
  - b) Value-regularized method
  - c) Implicit constraints/regularization
  - d) Model-based method
  - e) Conditional imitation method



# Policy constraint

- **Explicit policy constraint methods**

$$\sum_{(\mathbf{s}, \mathbf{a}, \mathbf{s}') \sim \mathcal{D}} \| Q(\mathbf{s}, \mathbf{a}) - (r(\mathbf{s}, \mathbf{a}) + \gamma \max_{\mathbf{a}'} Q(\mathbf{s}', \mathbf{a}')) \|^2$$

An intuitive idea: let  $\mathbf{a}'$  stay close to behavior policy

New objective:

$$\sum_{(\mathbf{s}, \mathbf{a}, \mathbf{s}') \sim \mathcal{D}} \| Q(\mathbf{s}, \mathbf{a}) - (r(\mathbf{s}, \mathbf{a}) + \gamma E_{\mathbf{a}' \sim \pi_{new}(\cdot | \mathbf{s}')} Q(\mathbf{s}', \mathbf{a}')) \|^2$$

$\pi_{new}(\mathbf{a} | \mathbf{s}) = \arg \max_{\pi} E_{\mathbf{a} \sim \pi(\mathbf{a} | \mathbf{s})} [Q(\mathbf{s}, \mathbf{a})]$  s.t.  $D_{KL}(\pi \| \pi_{\beta}) \leq \epsilon$



# Policy constraint

$$\pi_{new}(\mathbf{a} \mid \mathbf{s}) = \arg \max_{\pi} E_{\mathbf{a} \sim \pi(\mathbf{a} \mid \mathbf{s})} [Q(\mathbf{s}, \mathbf{a})] \quad \text{s.t. } D_{\text{KL}}(\pi \parallel \pi_{\beta}) \leq \epsilon$$

**Issue:** we don't know what  $\pi_{\beta}$  is

How do we implement constraints in RL?

- fit a policy to the data,
- *i.e.*, learn a proxy for  $\pi_{\beta}$  through imitation

Forms of policy constraints?

1. **support constraint:**  $\pi(\mathbf{a} \mid \mathbf{s}) > 0$  only if  $\pi_{\beta}(\mathbf{a} \mid \mathbf{s}) \geq \epsilon$   
+ close to what we want      - challenging to implement in practice
2. **KL divergence:**  $D_{KL}(\pi \parallel \pi_{\beta})$   
+ easy to implement      - not necessarily what we want



# Policy constraint

**Unconstrained optimization objective:**

$$\theta \leftarrow \arg \max_{\theta} E_{\mathbf{s} \sim D} [E_{\mathbf{a} \sim \pi_{\theta}(\mathbf{a}|\mathbf{s})} [Q(\mathbf{s}, \mathbf{a}) + \lambda \log \pi_{\beta}(\mathbf{a}|\mathbf{s})] + \lambda \mathcal{H}(\pi(\mathbf{a}|\mathbf{s}))]$$

Lagrange multiplier  $\lambda \geq 0$

- Since for the constraints we have

$$D_{\text{KL}}(\pi \| \pi_{\beta}) = E_{\pi} [\log \pi(\mathbf{a}|\mathbf{s}) - \log \pi_{\beta}(\mathbf{a}|\mathbf{s})] = -E_{\pi} [\log \pi_{\beta}(\mathbf{a}|\mathbf{s})] - \mathcal{H}(\pi)$$

**Modified reward function:**

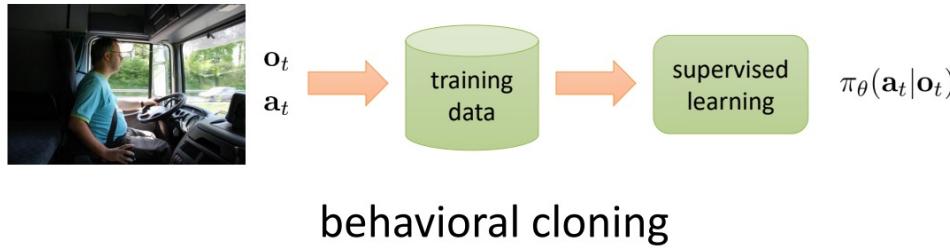
$$\bar{r}(s, a) = r(s, a) - D_{KL}(\pi, \pi_{\beta})$$

Penalize action distribution divergence



# Policy constraint

- Like imitation learning? No!

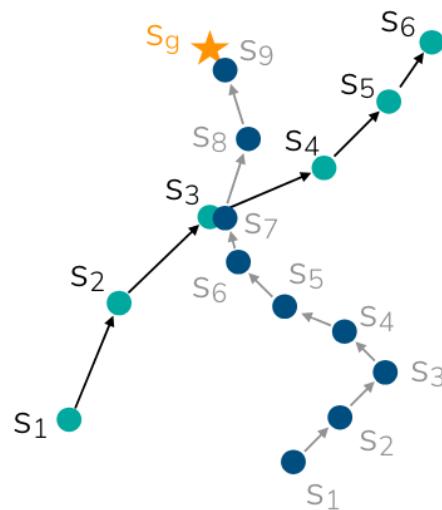


- Imitation methods can't outperform the expert
- Offline data may not be optimal!
- Offline RL **can leverage reward information** to outperform behavior policy.
- Good offline RL methods can ***stitch*** together good behaviors.



- Imitation learning vs. offline reinforcement learning

Good offline RL methods can ***stitch*** together good behaviors.



$s_1 \rightarrow s_3$  is good behavior

$s_7 \rightarrow s_9$  is good behavior

Offline RL methods can learn a policy that goes from  $s_1$  to  $s_9$ !



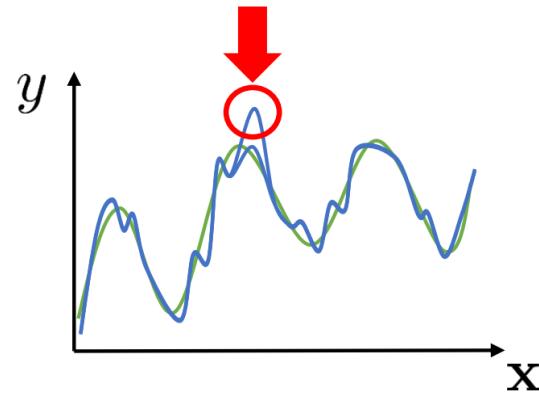
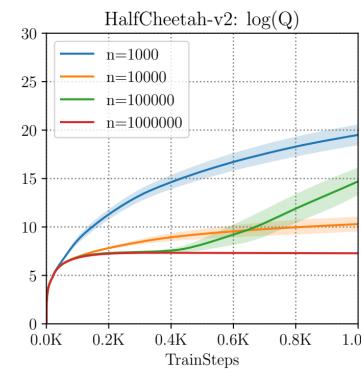
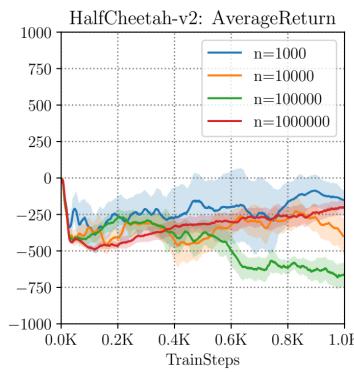
# Contents at a glance

1. Preliminary
2. Challenges of Offline RL
3. Solutions
  - a) Policy constraint methods
  - b) Value-regularized method
  - c) Implicit constraints/regularization
  - d) Model-based method
  - e) Conditional imitation method



# Conservative Q-learning (CQL)

- **Regularizing the Q-function directly**  
Avoid overestimation for out-of-distribution actions



General objective for conservative Q-function

$$\mathcal{E}(\mathcal{D}, Q_\phi) = \mathcal{E}(\mathcal{D}, Q_\phi) + \alpha \mathcal{C}(\mathcal{D}, Q_\phi)$$

Regular objective: Bellman error term  
Computed from dataset  $\mathcal{D}$

Conservative penalty term



# Conservative Q-learning (CQL)

- Regularizing the Q-function directly

$$\hat{Q}^\pi = \arg \min_Q \max_\mu E_{(s,a,s') \sim D} \left[ \underbrace{\left( Q(s, a) - (r(s, a) + \gamma E_\pi [Q(s', a')]) \right)^2}_{\text{standard critic update}} + \underbrace{\alpha E_{s \sim D, a \sim \mu(\cdot | s)} [Q(s, a)]}_{\text{push down on large Q-values}} \right]$$

- Push down large Q-values
- Can show that  $Q_\phi^\pi < Q^\pi$  (true Q-function) for large enough  $\alpha$



# Conservative Q-learning (CQL)

- Regularizing the Q-function directly

$$\hat{Q}^\pi = \arg \min_Q \max_\mu E_{(s,a,s') \sim D} \left[ \underbrace{\left( Q(s, a) - (r(s, a) + \gamma E_\pi[Q(s', a')]) \right)^2}_{\text{standard critic update}} \right] + \alpha E_{s \sim D, a \sim \mu(\cdot | s)} [Q(s, a)]$$

push down on large Q-values

$$- \underbrace{\alpha E_{(s,a) \sim D} [Q(s, a)]}_{\text{push up on Q-values for } (s, a) \text{ in the data}}$$

- No longer guaranteed that  $Q_\phi^\pi < Q^\pi$  for all  $(s, a)$
- Can show  $E_{a \sim \pi}[Q_\phi^\pi(s, a)] < E_{a \sim \pi}[Q^\pi(s, a)]$  for all  $s \in \mathcal{D}$



- Full algorithm of conservative Q-learning

- 
1. Update  $\hat{Q}^\pi$  using  $L_{CQL}$  using  $D$
  2. Update policy  $\pi$

If actions are discrete:  $\pi(\mathbf{a} \mid \mathbf{s}) = \begin{cases} 1 & \text{if } \mathbf{a} = \arg \max_{\bar{\mathbf{a}}} \hat{Q}(\mathbf{s}, \bar{\mathbf{a}}) \\ 0 & \text{otherwise} \end{cases}$

If actions are continuous:  $\theta \leftarrow \theta + \eta \nabla_\theta E_{\mathbf{s} \sim D, \mathbf{a} \sim \pi_\theta(\cdot | \mathbf{s})} [\hat{Q}(\mathbf{s}, \mathbf{a})]$



# Contents at a glance

1. Preliminary
2. Challenges of Offline RL
3. Solutions
  - a) Policy constraint methods
  - b) Value-regularized method
  - c) Implicit constraints/regularization
  - d) Model-based method
  - e) Conditional imitation method



# Policy Constraint

- **Problems in explicit policy constraints**

Behavior policy is unknown

Recall: learn a proxy for  $\pi_\beta$  through imitation

Estimation error for approximating behavior policy



**Not** to explicitly model behavior policy



# Implicit Policy Constraint

- Expected improvement in terms of advantage func.

$$\eta(\pi) = \mathbb{E}_{\mathbf{s} \sim d_\pi(\mathbf{s})} \mathbb{E}_{\mathbf{a} \sim \pi(\mathbf{a}|\mathbf{s})} [A^\mu(\mathbf{s}, \mathbf{a})] = \mathbb{E}_{\mathbf{s} \sim d_\pi(\mathbf{s})} \mathbb{E}_{\mathbf{a} \sim \pi(\mathbf{a}|\mathbf{s})} [\mathcal{R}_{\mathbf{s}, \mathbf{a}}^\mu - V^\mu(\mathbf{s})]$$

$$V^\mu(\mathbf{s}) = \int_a \boxed{\pi_\beta(\mathbf{a} \mid \mathbf{s}) \mathcal{R}_s^\mathbf{a}} d\mathbf{a}$$

Behavior policy

- Objective

$$\begin{aligned} \arg \max_{\pi} & \quad \int_{\mathbf{s}} d_\mu(\mathbf{s}) \int_{\mathbf{a}} \pi(\mathbf{a} \mid \mathbf{s}) [A^\mu] d\mathbf{a} d\mathbf{s} \\ \text{s.t.} & \quad \int_{\mathbf{s}} d_\mu(\mathbf{s}) D_{\text{KL}}(\pi(\cdot \mid \mathbf{s}) \parallel \pi_\beta(\cdot \mid \mathbf{s})) d\mathbf{s} \leq \epsilon. \end{aligned}$$



# Implicit Policy Constraint

- Lagrange function

$$\mathcal{L}(\pi, \beta) = \int_s d_\mu(s) \int_a (a | s) A^\mu da ds \pi + \beta (\epsilon - \int_s d_\mu(s) D_{KL}(\pi(\cdot | s) \| \pi_\beta(\cdot | s)) ds)$$



- Optimal value

$$\pi^*(a | s) = \frac{1}{Z(s)} \pi_\beta(a | s) \exp\left(\frac{1}{\lambda} A^\mu(s, a)\right)$$

hyperparameter



# Implicit Policy Constraint

- Learning a policy close to optimal

$$\pi^*(\mathbf{a} \mid \mathbf{s}) = \frac{1}{Z(\mathbf{s})} \pi_\beta(\mathbf{a} \mid \mathbf{s}) \exp\left(\frac{1}{\lambda} A^\mu(\mathbf{s}, \mathbf{a})\right)$$

$$\arg \min_{\pi} \mathbb{E}_{\mathbf{s} \sim \mathcal{D}} [\text{D}_{\text{KL}}(\pi^*(\cdot \mid \mathbf{s}) \parallel \pi(\cdot \mid \mathbf{s}))]$$



$$\hat{\pi} \leftarrow \arg \max_{\pi} E_{\mathbf{s}, \mathbf{a} \sim D} [\log \pi(\mathbf{a} \mid \mathbf{s}) \exp\left(\frac{1}{\lambda} A^\mu\right)]$$

hyperparameter

$$D_{\text{KL}}(\pi \parallel \pi_\beta) = E_\pi [\log \pi(\mathbf{a} \mid \mathbf{s}) - \log \pi_\beta(\mathbf{a} \mid \mathbf{s})] = -E_\pi [\log \pi_\beta(\mathbf{a} \mid \mathbf{s})] - \mathcal{H}(\pi)$$



- **Advantage-weighted regression**

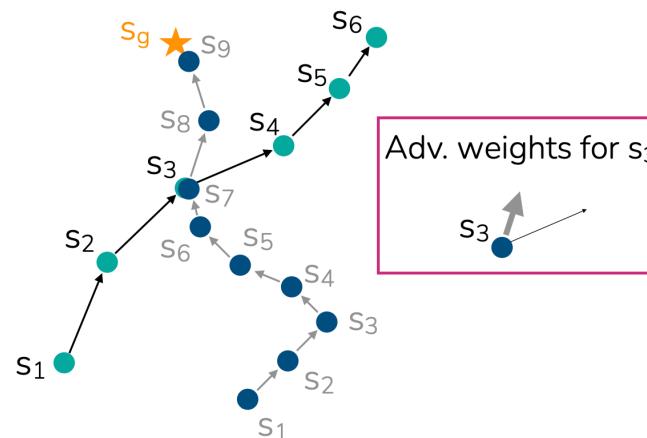
Another perspective

$$\hat{\pi} \leftarrow \arg \max_{\pi} E_{\mathbf{s}, \mathbf{a} \sim D} [\log \pi(\mathbf{a} | \mathbf{s}) \exp(\frac{1}{\lambda} A^\mu)]$$



$$\theta \leftarrow \arg \max_{\theta} E_{\mathbf{s}, \mathbf{a} \sim D} [\log \pi_\theta(\mathbf{a} | \mathbf{s}) \exp(A(\mathbf{s}, \mathbf{a}))]$$

standard imitation learning    with advantage weights





- How to estimate the advantage function

Estimate  $V^{\pi_\beta}(s)$  with Monte Carlo,  $\min_V E_{(\mathbf{s}, \mathbf{a}) \sim D} \left[ (R_{\mathbf{s}, \mathbf{a}} - V(\mathbf{s}))^2 \right]$

Approximate  $\hat{A}^{\pi_\beta}(\mathbf{s}, \mathbf{a}) = R_{\mathbf{s}, \mathbf{a}} - V(\mathbf{s})$



- + Simple
- + Avoids querying or training on any OOD actions!

- Monte Carlo estimation is noisy
- $\hat{A}^{\pi_\beta}$  assumes weaker policy than  $\hat{A}^{\pi_\theta}$



- **Advantage weighted actor critic**

Estimate advantage function with TD updates instead of Monte Carlo?

1. Estimate  $Q^\pi$ -function:  $\min_Q E_{(s,a,s') \sim D} \left[ \left( Q(s, a) - \left( r + \gamma E_{a' \sim \pi(\cdot|s)} [Q(s', a')] \right) \right)^2 \right]$
2. Estimate advantage as:  $\hat{A}^\pi(s, a) = \hat{Q}^\pi(s, a) - E_{\bar{a} \sim \pi(\cdot|s)} [\hat{Q}^\pi(s, \bar{a})]$
3. Update policy as before:  $\hat{\pi} \leftarrow \arg \max_\pi E_{s,a \sim D} \left[ \log \pi(a | s) \exp \left( \frac{1}{\alpha} \hat{A}^\pi(s, a) \right) \right]$

- + Policy still only trained on actions in data.
- + Temporal difference updates instead of Monte Carlo.
- What might go wrong?
  - Possibly querying OOD actions!



- Can we estimate Q without querying OOD actions?

**AWAC:** Estimate  $Q$ -function:  $\min_Q E_{(s,a,s') \sim D} \left[ \left( Q(s, a) - \left( r + \gamma E_{\substack{a' \sim \pi(\cdot|s) \\ a' \sim D}} [Q(s', a')] \right) \right)^2 \right]$

“SARSA algorithm”

*Evaluate the Q-function only on the  $(s, a)$  in the offline dataset*



The core idea of **Implicit Q-Learning**

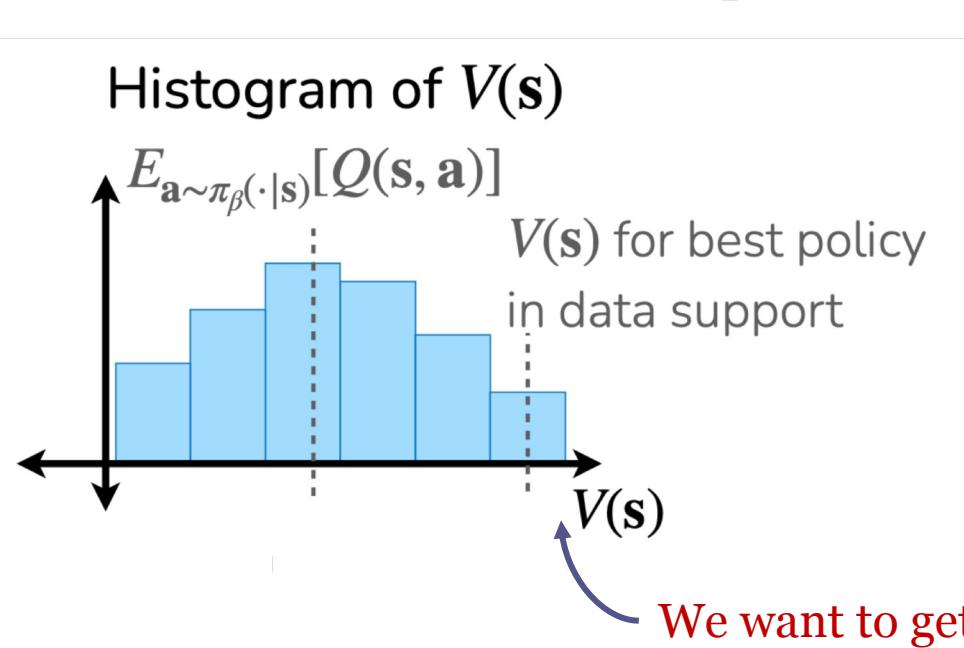


- Can we estimate  $Q$  for a policy  $\pi'(a | s)$  that is better than  $\pi_\beta(a | s)$  ?

We don't need to explicitly modeling  $\pi'(a | s)$

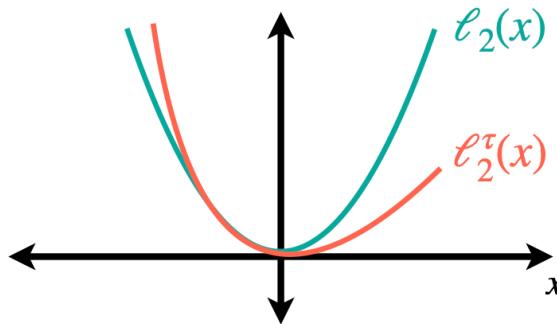
SARSA update:  $\hat{Q}^{\pi_\beta} \leftarrow \arg \min_Q E_{(s,a,s',a') \sim D} \left[ \left( Q(s, a) - (r + \gamma \underline{Q(s', a')}) \right)^2 \right]$

a sample of  $V^{\pi_\beta}(s')$

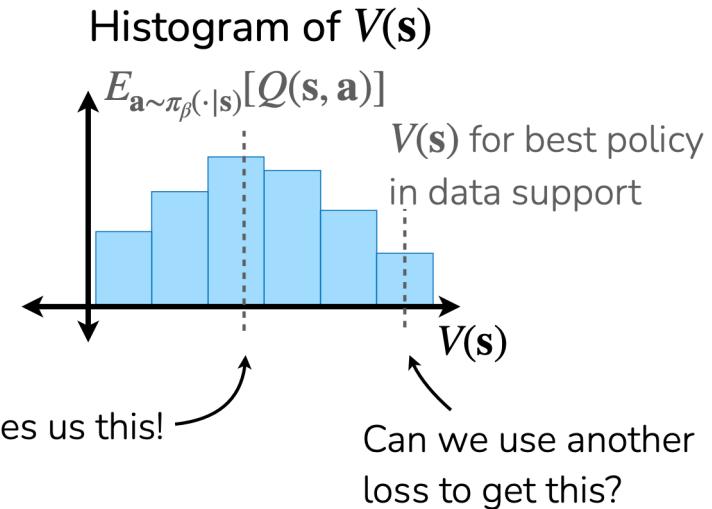




- Key idea: Expectile regression

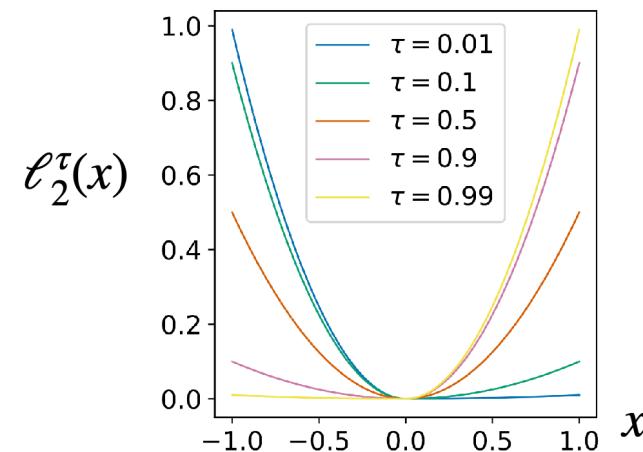


$\ell_2$  loss gives us this!



Expectile regression loss:

$$\ell_2^\tau(x) = \begin{cases} (1 - \tau)x^2 & \text{if } x < 0 \\ \tau x^2 & \text{otherwise} \end{cases}$$





- **Implicit Q-Learning**

$$L(\theta) = \mathbb{E}_{(s, a, s', a') \sim \mathcal{D}} \left[ l_2^\tau \left( r(s, a) + \gamma Q_{\hat{\theta}}(s', a') - Q_\theta(s, a) \right) \right]$$

Expectile regression loss:

$$\ell_2^\tau(x) = \begin{cases} (1 - \tau)x^2 & \text{if } x < 0 \\ \tau x^2 & \text{otherwise} \end{cases}$$

A large  $Q_\theta(s, a)$   a good action  $a$

a “lucky” sample that happened to have transitioned into a good state.



$$L_V(\psi) = \mathbb{E}_{(s, a) \sim \mathcal{D}} [L_2^\tau(Q_{\hat{\theta}}(s, a) - V_\psi(s))].$$

$$L_Q(\theta) = \mathbb{E}_{(s, a, s') \sim \mathcal{D}} [(r(s, a) + \gamma V_\psi(s') - Q_\theta(s, a))^2].$$



## • Implicit Q-Learning

### Full algorithm

Fit  $V$  with expectile loss:  $\hat{V}(\mathbf{s}) \leftarrow \arg \min_V E_{(\mathbf{s}, \mathbf{a}) \sim D} \left[ \ell_2^\tau \left( V(\mathbf{s}) - \hat{Q}(\mathbf{s}, \mathbf{a}) \right) \right]$  using small  $\tau < 0.5$

Update  $Q$  with typical MSE loss:  $\hat{Q}(\mathbf{s}, \mathbf{a}) \leftarrow \arg \min_Q E_{(\mathbf{s}, \mathbf{a}, \mathbf{s}') \sim D} \left[ \left( Q(\mathbf{s}, \mathbf{a}) - (r + \gamma \hat{V}(\mathbf{s}')) \right)^2 \right]$

Extract policy with AWR:  $\hat{\pi} \leftarrow \arg \max_\pi E_{\mathbf{s}, \mathbf{a} \sim D} \left[ \log \pi(\mathbf{a} | \mathbf{s}) \exp \left( \frac{1}{\alpha} \left( \hat{Q}(\mathbf{s}, \mathbf{a}) - \hat{V}(\mathbf{s}) \right) \right) \right]$

- + Never need to query OOD actions!
- + Policy (still) only trained on actions in data.
- + Decoupling actor & critic training —> computationally fast

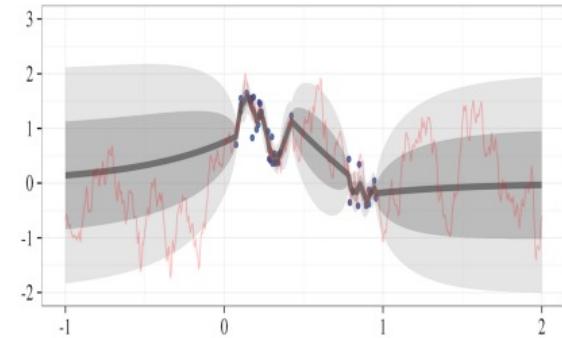
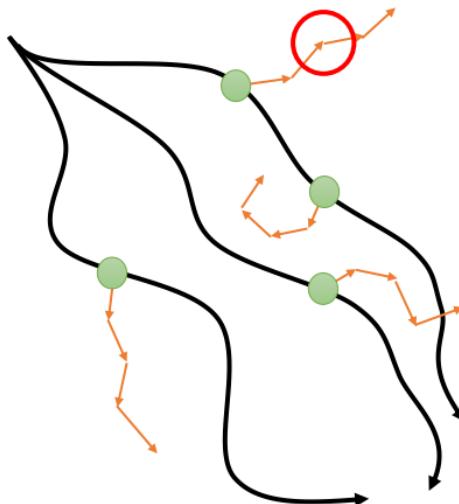


# Contents at a glance

1. Preliminary
2. Challenges of Offline RL
3. Solutions
  - a) Policy constraint methods
  - b) Value-regularized method
  - c) Implicit constraints/regularization
  - d) **Model-based method**
  - e) Conditional imitation method

# Model-Based Offline RL

- **Advantage of model-based RL under offline setting**
  - Supervised learning is more stable than bootstrapping RL training
  - Learned environment model can provide uncertainty estimation

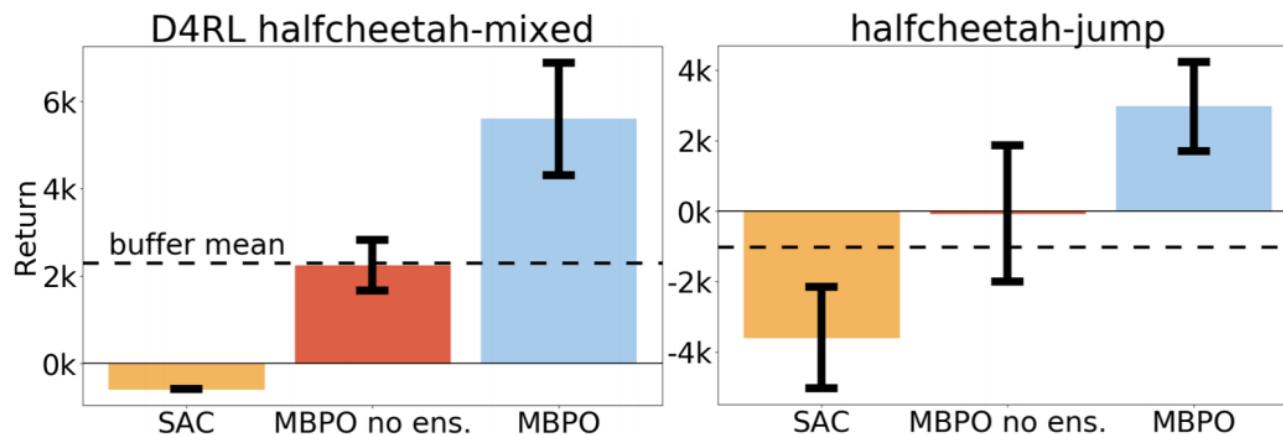


Gaussian process regression

Query of OOD states can be measured by uncertainty

# Model-Based Offline RL

- Model-based RL is well-suited for offline RL




---

**Algorithm 1** Framework for Model-based Offline Policy Optimization (MOPO) with Reward Penalty

**Require:** Dynamics model  $\hat{T}$  with admissible error estimator  $u(s, a)$ ; constant  $\lambda$ .

- 1: Define  $\tilde{r}(s, a) = r(s, a) - \lambda u(s, a)$ . Let  $\tilde{M}$  be the MDP with dynamics  $\hat{T}$  and reward  $\tilde{r}$ .
  - 2: Run any RL algorithm on  $\tilde{M}$  until convergence to obtain  $\hat{\pi} = \operatorname{argmax}_{\pi} \eta_{\tilde{M}}(\pi)$
-

## ■ MOPO: Reward penalty when model is uncertain

- **Key idea:** build a lower bound for the expected return of a policy  $\pi$  under the true dynamics and then maximize the lower bound over  $\pi$

**Lemma 4.1** (Telescoping lemma). *Let  $M$  and  $\widehat{M}$  be two MDPs with the same reward function  $r$ , but different dynamics  $T$  and  $\widehat{T}$  respectively. Let  $G_{\widehat{M}}^\pi(s, a) := \mathbb{E}_{s' \sim \widehat{T}(s, a)}[V_M^\pi(s')] - \mathbb{E}_{s' \sim T(s, a)}[V_M^\pi(s')]$ . Then,*

$$\eta_{\widehat{M}}(\pi) - \eta_M(\pi) = \gamma \mathbb{E}_{(s, a) \sim \rho_{\widehat{T}}^\pi} \left[ G_{\widehat{M}}^\pi(s, a) \right] \quad (1)$$

As an immediate corollary, we have

$$\eta_M(\pi) = \mathbb{E}_{(s, a) \sim \rho_{\widehat{T}}^\pi} \left[ r(s, a) - \gamma G_{\widehat{M}}^\pi(s, a) \right] \geq \mathbb{E}_{(s, a) \sim \rho_{\widehat{T}}^\pi} \left[ r(s, a) - \gamma |G_{\widehat{M}}^\pi(s, a)| \right] \quad (2)$$

- Measure the difference between two MDPs

$$|G_{\widehat{M}}^\pi(s, a)| \leq \sup_{f \in \mathcal{F}} \left| \mathbb{E}_{s' \sim \widehat{T}(s, a)}[f(s')] - \mathbb{E}_{s' \sim T(s, a)}[f(s')] \right| =: d_{\mathcal{F}}(\widehat{T}(s, a), T(s, a)),$$

$$\eta_M(\pi) \geq \mathbb{E}_{(s, a) \sim \rho_{\widehat{T}}^\pi} \left[ r(s, a) - \gamma |G_{\widehat{M}}^\pi(s, a)| \right] \geq \mathbb{E}_{(s, a) \sim \rho_{\widehat{T}}^\pi} [r(s, a) - \lambda u(s, a)]$$

where  $d_{\mathcal{F}}(\widehat{T}(s, a), T(s, a)) \leq u(s, a)$

- T. Yu, T. Ma, et.al. MOPO: Model-based offline policy optimization. NeurIPS, 2020.

## ■ COMBO: Conservative Q value w.r.t OOD states and actions

- Model-free conservative Q-learning

➤ Policy Evaluation

$$Q^{k+1} \leftarrow \arg \min_Q \beta (\mathbb{E}_{s \sim \mathcal{D}, a \sim \mu(\cdot|s)} [Q(s, a)] - \mathbb{E}_{s, a \sim \mathcal{D}} [Q(s, a)]) \quad \text{➤ Policy improvement}$$

$$+ \frac{1}{2} \mathbb{E}_{s, a, s' \sim \mathcal{D}} \left[ (Q(s, a) - \hat{\mathcal{B}}^\pi Q^k(s, a))^2 \right], \quad (3) \quad \pi \leftarrow \arg \max_{\pi'} \mathbb{E}_{s \sim \mathcal{D}, a \sim \pi'(\cdot|s)} [\hat{Q}^\pi(s, a)]$$

$\mathcal{D}$  is offline dataset,  $\mu$  is a sampling distribution:  $\mu(a, s) = d^{\pi_\beta}(s)\mu(a|s)$

- Model-based

We have  $\mathbb{E}_{\pi(a|s)}[\hat{Q}^\pi(s, a)] \leq V^\pi(s)$  when  $\mu(a|s) = \pi(a|s)$

➤ Policy Evaluation

$$\hat{Q}^{k+1} \leftarrow \arg \min_Q \beta (\mathbb{E}_{s, a \sim \rho(s, a)} [Q(s, a)] - \mathbb{E}_{s, a \sim \mathcal{D}} [Q(s, a)]) \quad \rho(s, a) = d_{\hat{\mathcal{M}}}^\pi(s)\pi(a|s)$$

$$+ \frac{1}{2} \mathbb{E}_{s, a, s' \sim d_f} \left[ (Q(s, a) - \hat{\mathcal{B}}^\pi \hat{Q}^k(s, a))^2 \right]. \quad (4) \quad d_f^\mu(s, a) := f d(s, a) + (1 - f) d_{\hat{\mathcal{M}}}^\mu(s, a),$$

$\rho$  is a sampling distribution over learned model  $\hat{\mathcal{M}}$ .

- Tianhe Yu, and Chelsea Finn et al. Combo: Conservative offline model-based policy optimization. arXiv preprint arXiv:2102.08363, 2021.



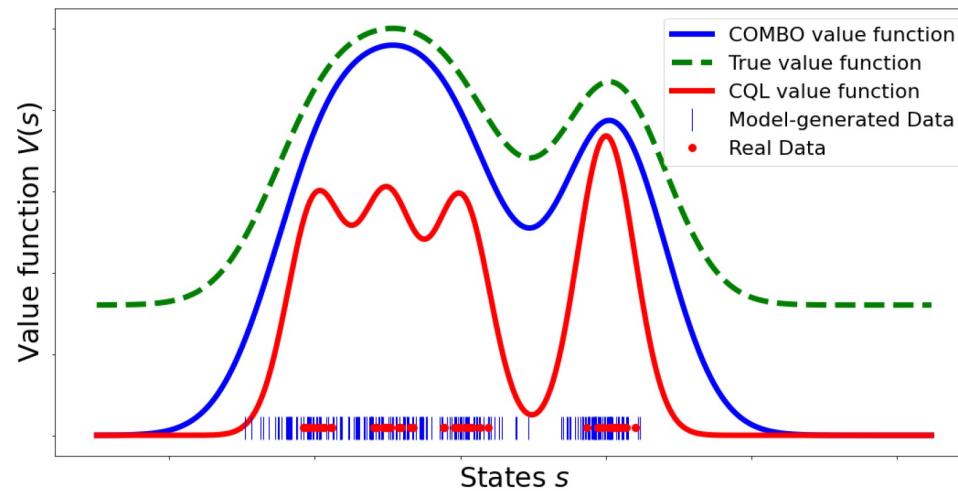
## ■ COMBO: Conservative Q value w.r.t OOD states and actions

### ➤ Policy Evaluation

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

$$d_f^\mu(\mathbf{s}, \mathbf{a}) := f d(\mathbf{s}, \mathbf{a}) + (1 - f) d_{\hat{\mathcal{M}}}^\mu(\mathbf{s}, \mathbf{a}),$$

$\rho$  is a sampling distribution over learned model  $\hat{\mathcal{M}}$ .



- Aviral Kumar, and Sergey Levine, et al.. Conservative q-learning for offline reinforcement learning. arXiv preprint arXiv:2006.04779, 2020.
- Tianhe Yu, and Chelsea Finn et al. Combo: Conservative offline model-based policy optimization. arXiv preprint arXiv:2102.08363, 2021.



# Contents at a glance

1. Preliminary
2. Challenges of Offline RL
3. Solutions
  - a) Policy constraint methods
  - b) Value-regularized method
  - c) Implicit constraints/regularization
  - d) Model-based method
  - e) Conditional imitation method



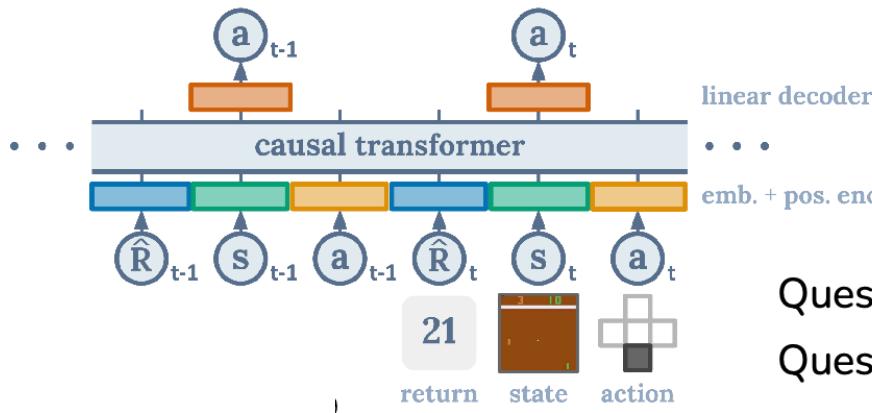
# Decision Transformer

1. Imitate entire dataset:  $\max_{\pi} \sum_{(s,a) \in D} \log \pi(a | s, \underline{R_{s,a}})$

Condition policy on (empirical) return to go.

- Policy will learn to mimic **good** and **poor** behaviors (and everything in between!)
- Pass in high return at test time
- Can use a sequence model:

Referred to as: upside-down RL, reward-conditioned policies, decision transformers



Question: Can this approach do data stitching?

Question: When would a sequence model be helpful?



# Diffusion model for Offline RL

- **Diffusion model**

Forward process, add noise: original image → Gaussian noise

$$q(\mathbf{x}_{k+1} | \mathbf{x}_k) := \mathcal{N}(\mathbf{x}_{k+1}; \sqrt{\alpha_k} \mathbf{x}_k, (1 - \alpha_k) \mathbf{I})$$

Reverse process, denoise (generation): Gaussian noise → image

$$p_{\theta}(\mathbf{x}_{k-1} | \mathbf{x}_k) := \mathcal{N}(\mathbf{x}_{k-1} | \underline{\mu_{\theta}(\mathbf{x}_k, k)}, \Sigma_k)$$

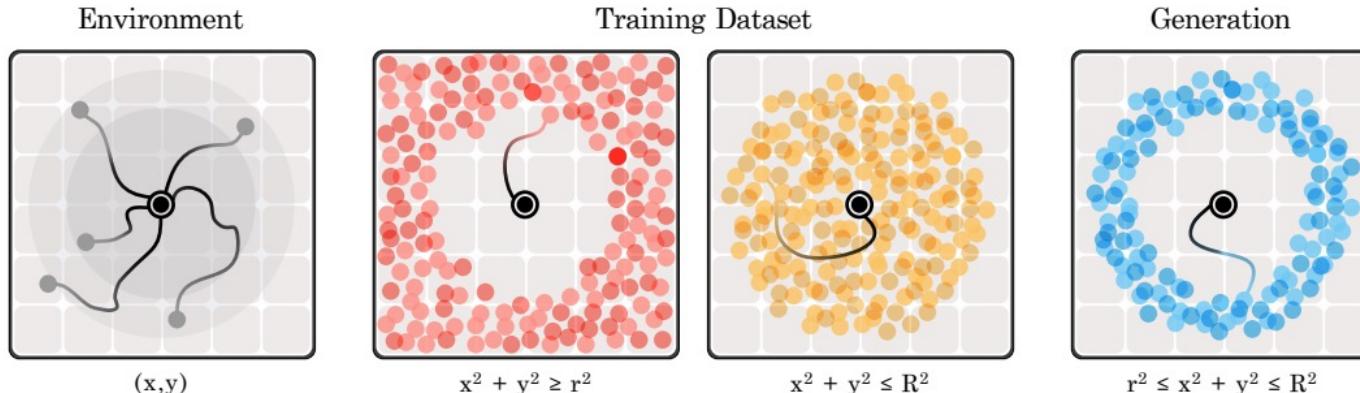
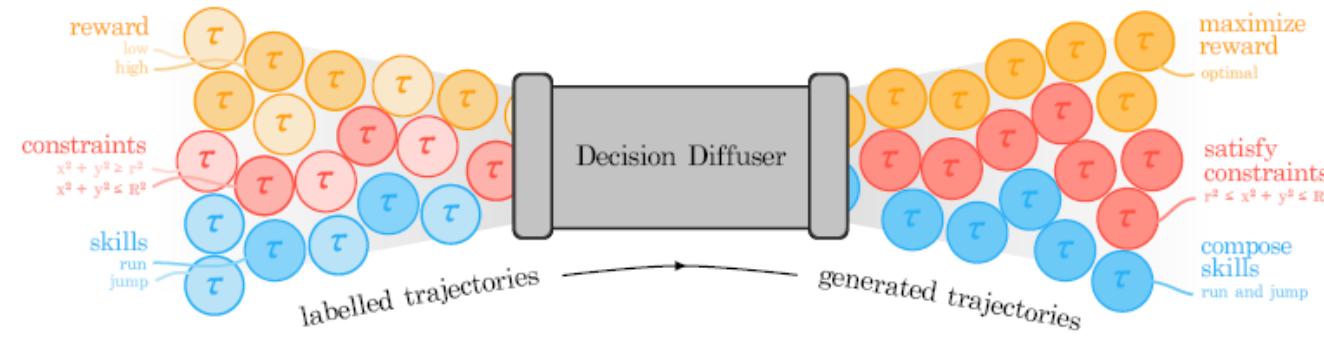
Guided diffusion: modeling distribution  $q(x|y)$  for generating samples with attributes of label  $y$

[1] Ho J, Jain A, Abbeel P. Denoising diffusion probabilistic models[J]. Advances in neural information processing systems, 2020, 33: 6840-6851.



# Diffusion model for Offline RL

- Decision Making using Conditional Generative Modeling





## Q & A



Many Thanks