

Offline Reinforcement Learning

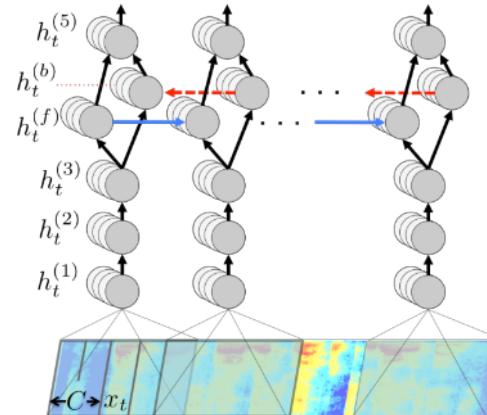
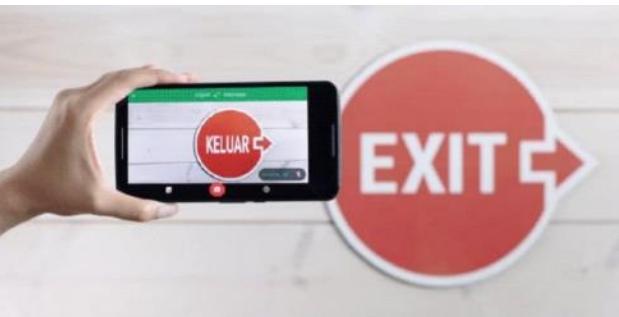
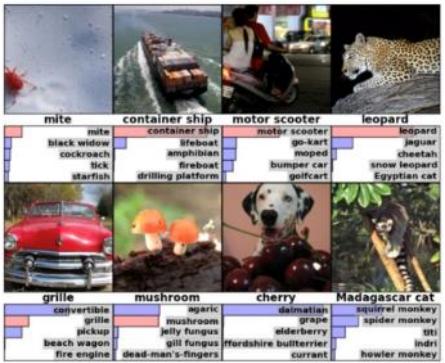
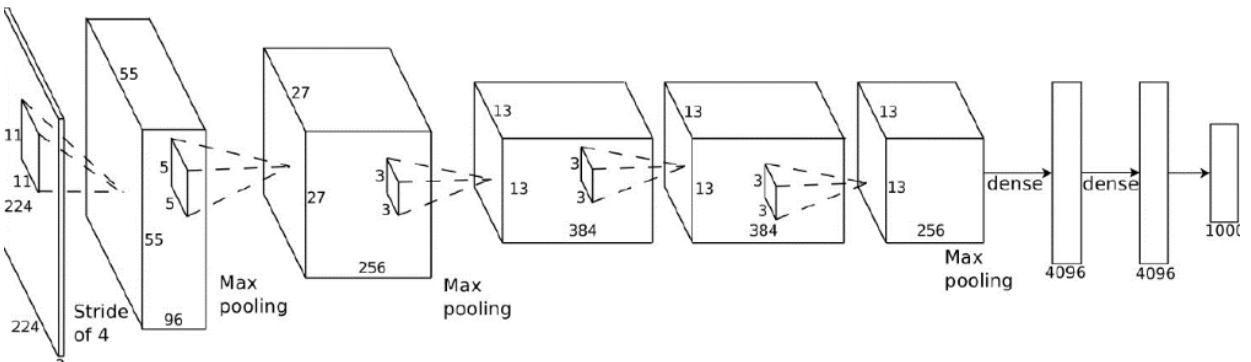
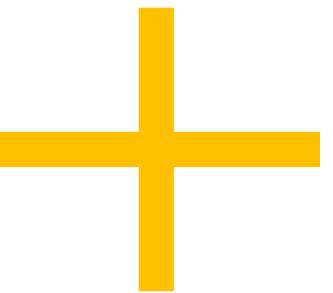
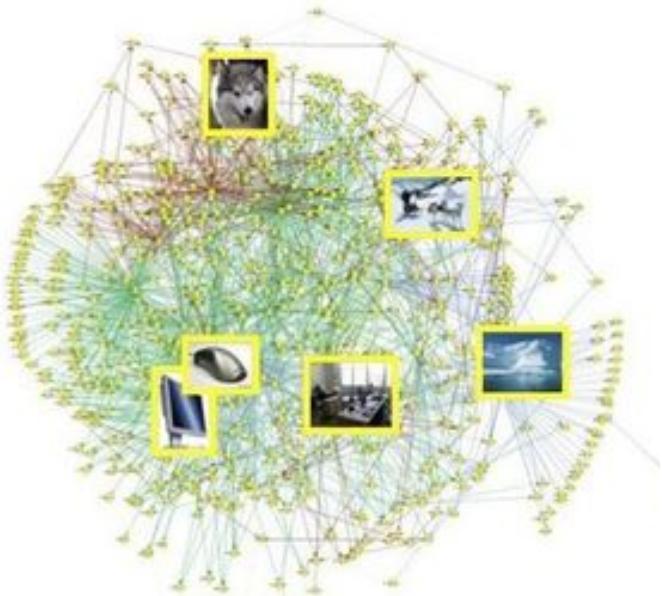
From Algorithms to Practical Challenges

NeurIPS 2020 Tutorial

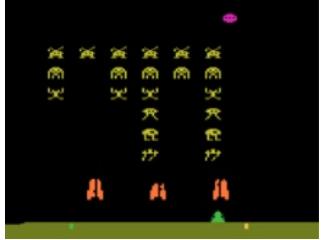
Aviral Kumar Sergey Levine
UC Berkeley



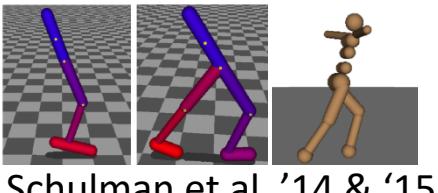
What makes modern machine learning work?



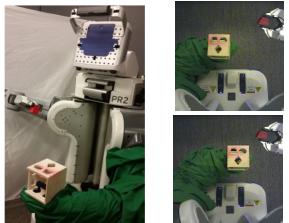
What about reinforcement learning?



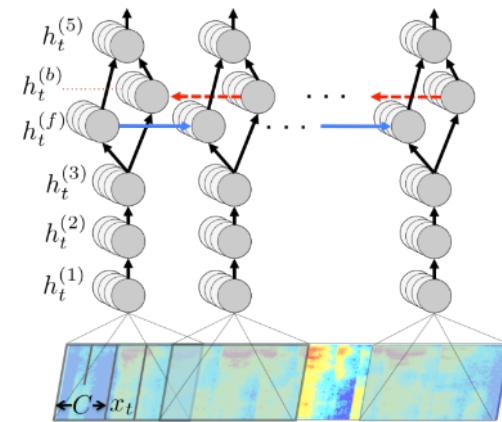
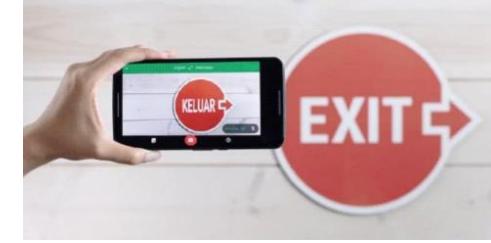
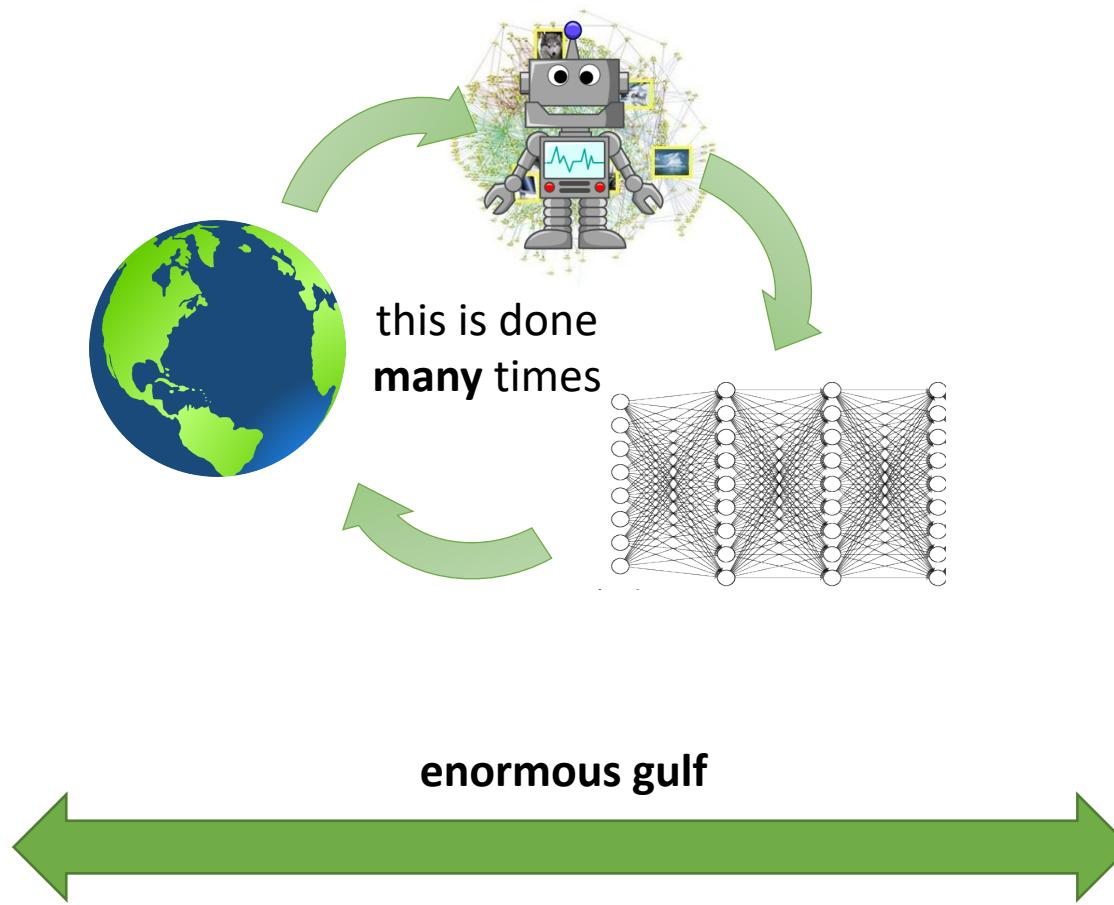
Mnih et al. '13



Schulman et al. '14 & '15



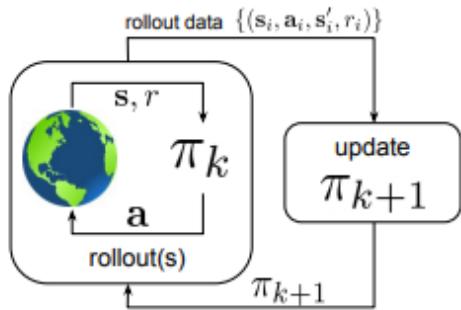
Levine*, Finn*, et al. '16



Can we develop data-driven RL methods?

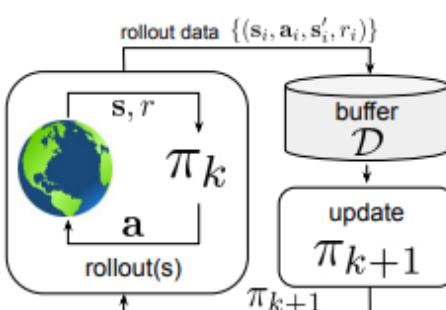
policy interacts with env
and updates the policy,
then collects new data

on-policy RL

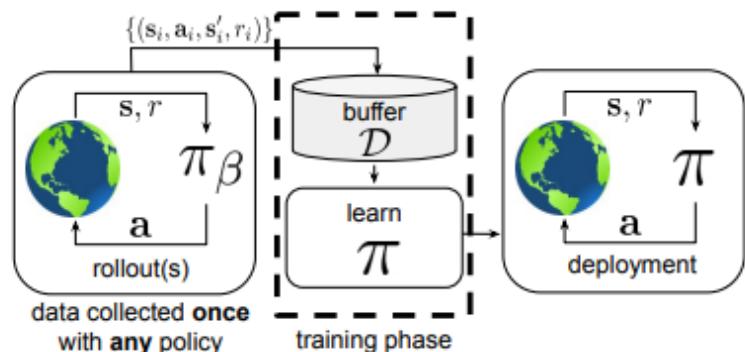


can use data collected by
other (usually previous)
policies. Still interacts
with the world.

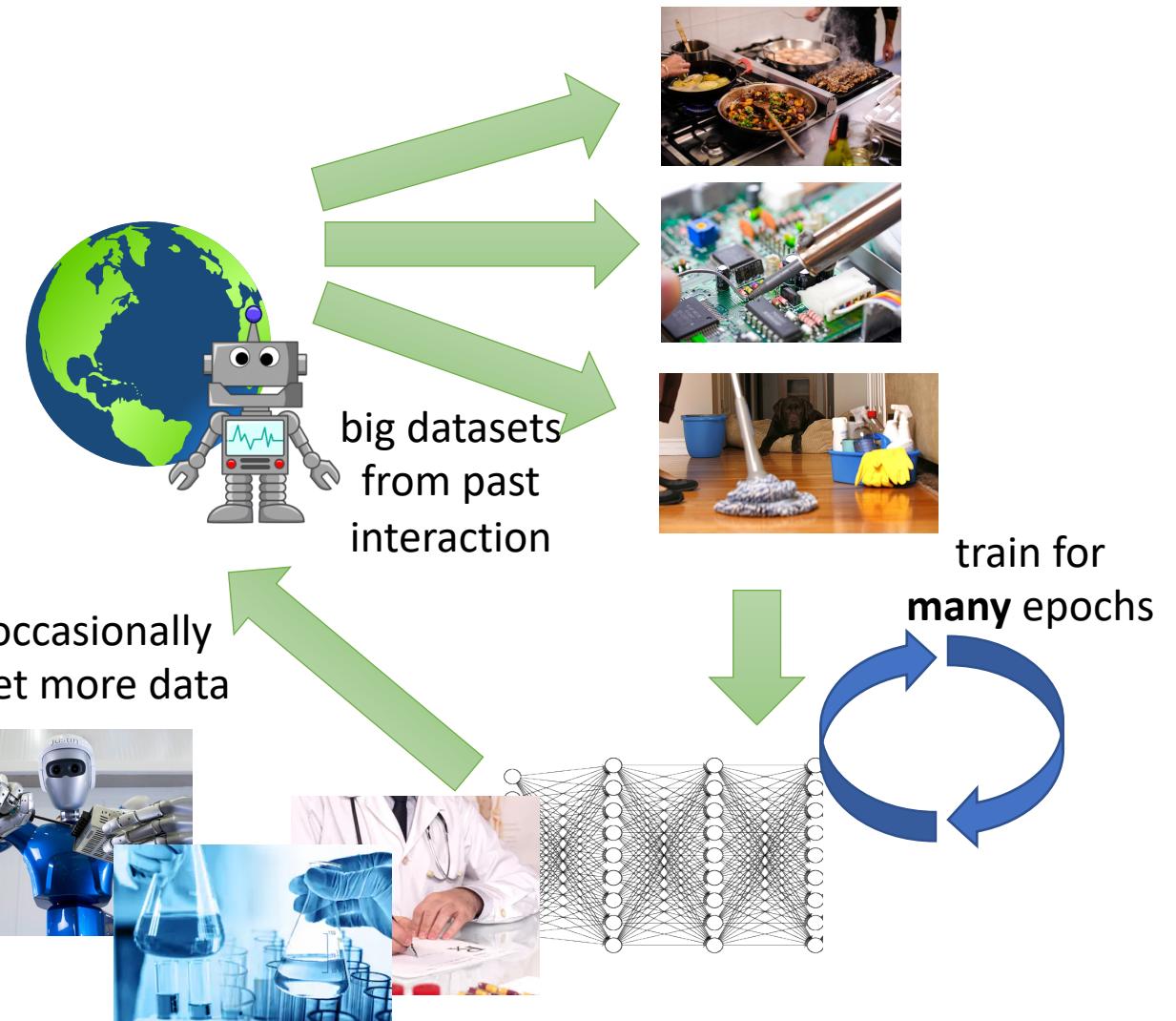
off-policy RL



offline reinforcement learning



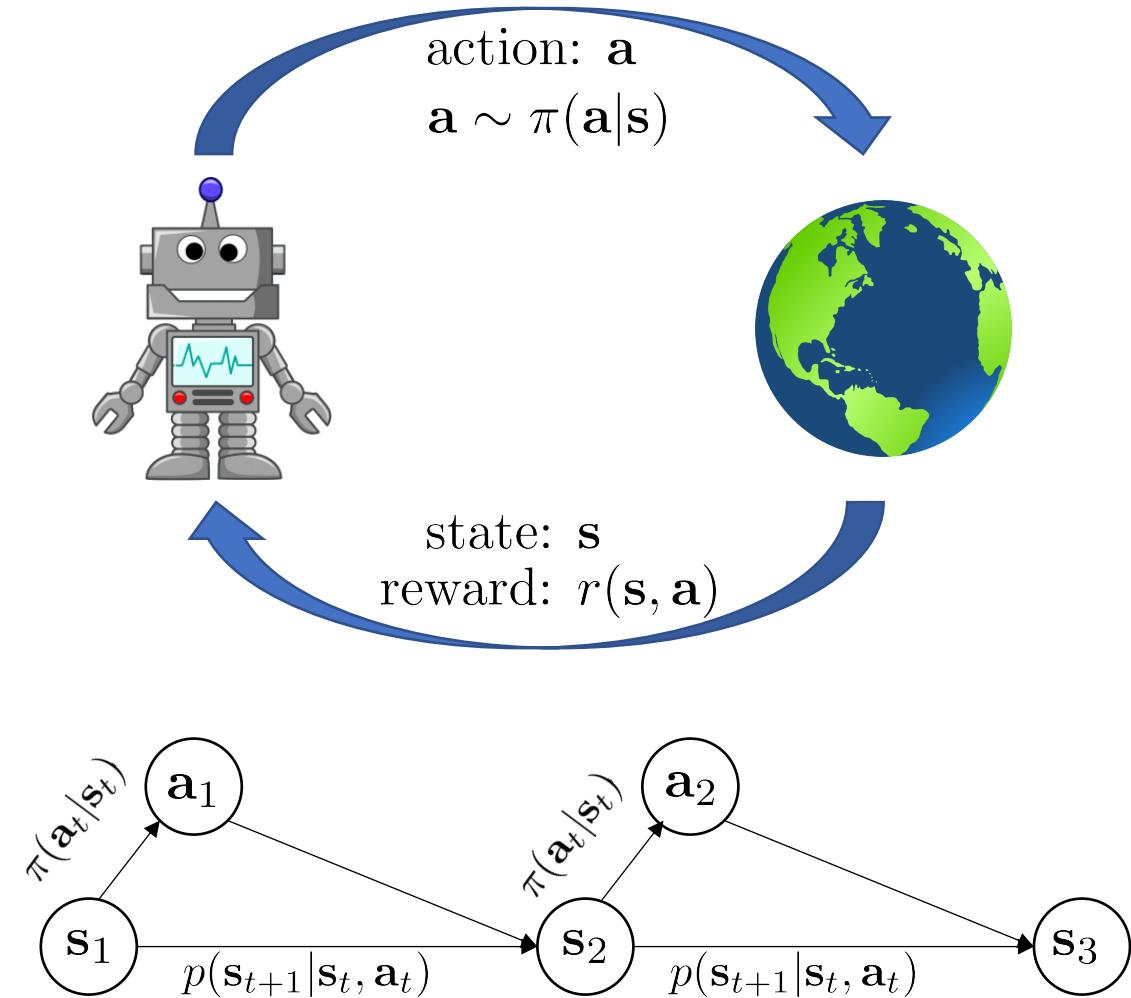
new data is not added. Behavior policy adds data to the buffer. We don't know what the policy is



1. Reinforcement learning primer
2. The offline reinforcement learning problem
3. Classic batch reinforcement learning algorithms
 - i. Importance-sampled policy gradients
 - ii. Least-squares temporal difference (LSTD)
4. Modern offline reinforcement learning algorithms
 - i. Policy constraint methods
 - ii. Model-based offline RL
 - iii. Value function regularization
 - iv. Uncertainty-based methods
 - v. Comparisons and discussion
5. Off-policy evaluation and model selection
 - i. Problem statement
 - ii. Direct methods and importance sampling methods
 - iii. Doubly robust methods
6. How should we evaluate offline RL methods?
 - i. When and why does offline RL work?
 - ii. How do we benchmark offline RL methods?
7. Conclusions, discussion, and open problems

- 1. Reinforcement learning primer**
2. The offline reinforcement learning problem
3. Classic batch reinforcement learning algorithms
 - i. Importance-sampled policy gradients
 - ii. Least-squares temporal difference (LSTD)
4. Modern offline reinforcement learning algorithms
 - i. Policy constraint methods
 - ii. Model-based offline RL
 - iii. Value function regularization
 - iv. Uncertainty-based methods
 - v. Comparisons and discussion
5. Off-policy evaluation and model selection
 - i. Problem statement
 - ii. Direct methods and importance sampling methods
 - iii. Doubly robust methods
6. How should we evaluate offline RL methods?
 - i. When and why does offline RL work?
 - ii. How do we benchmark offline RL methods?
7. Conclusions, discussion, and open problems

Policies and objectives



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

Annotations for the RL objective equation:

- could be ∞ : arrow points to the summand $E_{s_t \sim d^{\pi}, a_t \sim \pi(a|s)} [\gamma^t r(s_t, a_t)]$.
- discount factor: arrow points to the discount term γ^t .
- state distribution under π : arrow points to the expectation operator $E_{s_t \sim d^{\pi}, a_t \sim \pi(a|s)}$.

some definitions:

$s \in \mathcal{S}$ – discrete or continuous state

$a \in \mathcal{A}$ – discrete or continuous action

$\tau = \{s_0, a_0, s_1, a_1, \dots, s_T, a_T\}$ – trajectory

$$\pi(s_0, a_0, \dots, s_T, a_T) = p(s_1) \underbrace{\prod_{t=0}^T \pi(a_t|s_t)p(s_{t+1}|s_t, a_t)}_{\pi(\tau)}$$

$d_t^{\pi}(s_t)$ – state marginal of $\pi(\tau)$ at t

$d^{\pi}(s) = \frac{1}{1-\gamma} \sum_{t=0}^T \gamma^t d_t^{\pi}(s_t)$ – “visitation frequency”

Direct policy gradient

RL objective: $\max_{\pi} \sum_{t=0}^T E_{\mathbf{s}_t \sim d^\pi(\mathbf{s}), \mathbf{a}_t \sim \pi(\mathbf{a}|\mathbf{s})} [\gamma^t r(\mathbf{s}_t, \mathbf{a}_t)]$

↑ exactly the same thing!

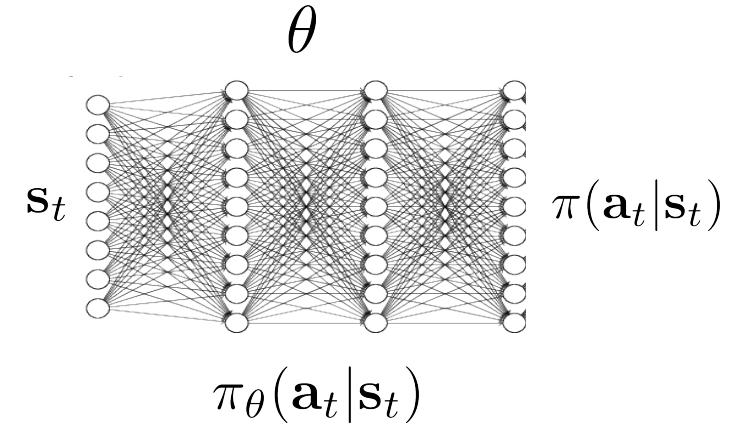
$$J(\theta) = E_{\tau \sim \pi_\theta(\tau)} \left[\sum_{t=0}^T \gamma^t r(\mathbf{s}_t, \mathbf{a}_t) \right] \approx \frac{1}{N} \sum_{i=1}^N \sum_{t=0}^T \gamma^t r(\mathbf{s}_{t,i}, \mathbf{a}_{t,i})$$

$$\nabla_{\theta} J(\theta) = E_{\tau \sim \pi_\theta(\tau)} \left[\nabla_{\theta} \log \pi_{\theta}(\tau) \sum_{t=0}^T \gamma^t r(\mathbf{s}_t, \mathbf{a}_t) \right] \text{ simple algebraic derivation}$$

$$= E_{\tau \sim \pi_\theta(\tau)} \left[\left(\sum_{t=0}^T \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_t | \mathbf{s}_t) \right) \left(\sum_{t=0}^T \gamma^t r(\mathbf{s}_t, \mathbf{a}_t) \right) \right] \text{ from definition of } \tau$$

$$= E_{\tau \sim \pi_\theta(\tau)} \left[\sum_{t=0}^T \nabla_{\theta} \gamma^t \log \pi_{\theta}(\mathbf{a}_t | \mathbf{s}_t) \left(\sum_{t'=t}^T \gamma^{t'-t} r(\mathbf{s}_t, \mathbf{a}_t) \right) \right]$$

$$= E_{\tau \sim \pi_\theta(\tau)} \left[\sum_{t=0}^T \nabla_{\theta} \gamma^t \log \pi_{\theta}(\mathbf{a}_t | \mathbf{s}_t) \hat{Q}(\mathbf{s}_t, \mathbf{a}_t) \right]$$



Policy gradient in practice

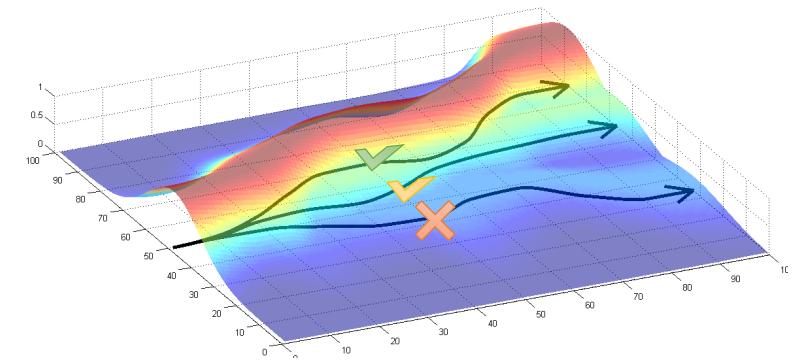
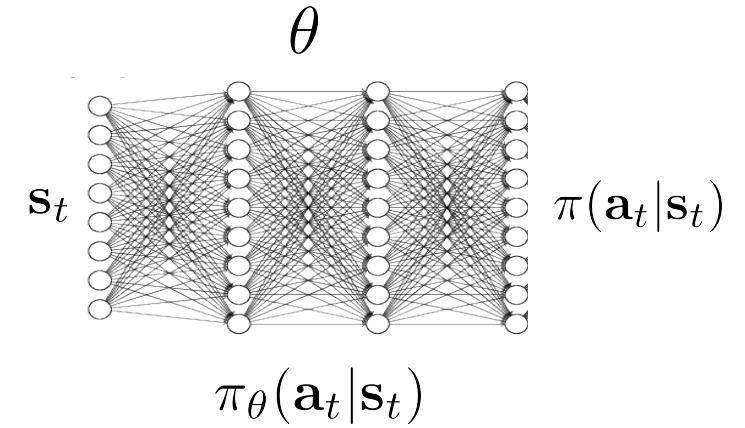
RL objective: $\max_{\pi} \sum_{t=0}^T E_{\mathbf{s}_t \sim d^\pi(\mathbf{s}), \mathbf{a}_t \sim \pi(\mathbf{a}|\mathbf{s})} [\gamma^t r(\mathbf{s}_t, \mathbf{a}_t)]$

$$\nabla_{\theta} J(\theta) = E_{\tau \sim \pi_{\theta}(\tau)} \left[\sum_{t=0}^T \nabla_{\theta} \gamma^t \log \pi_{\theta}(\mathbf{a}_t | \mathbf{s}_t) \hat{Q}(\mathbf{s}_t, \mathbf{a}_t) \right]$$

$$\approx \frac{1}{N} \sum_{i=1}^N \sum_{t=0}^T \nabla_{\theta} \gamma^t \log \pi_{\theta}(\mathbf{a}_{t,i} | \mathbf{s}_{t,i}) \hat{Q}(\mathbf{s}_{t,i}, \mathbf{a}_{t,i})$$

REINFORCE algorithm:

- 1. sample $\{\tau^i\}$ from $\pi_{\theta}(\mathbf{a}_t | \mathbf{s}_t)$ (run the policy)
- 2. $\nabla_{\theta} J(\theta) \approx \sum_i (\sum_t \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_t^i | \mathbf{s}_t^i)) (\sum_t r(\mathbf{s}_t^i, \mathbf{a}_t^i))$
- 3. $\theta \leftarrow \theta + \alpha \nabla_{\theta} J(\theta)$



Value functions and Q-functions

$$\nabla_{\theta} J(\theta) = E_{\tau \sim \pi_{\theta}(\tau)} \left[\sum_{t=0}^T \nabla_{\theta} \gamma^t \log \pi_{\theta}(\mathbf{a}_t | \mathbf{s}_t) \hat{Q}(\mathbf{s}_t, \mathbf{a}_t) \right] = \frac{1}{1-\gamma} E_{\mathbf{s} \sim d^{\pi_{\theta}}(\mathbf{s}), \mathbf{a} \sim \pi_{\theta}(\mathbf{a} | \mathbf{s})} [\nabla_{\theta} \log \pi_{\theta}(\mathbf{a} | \mathbf{s}) Q^{\pi_{\theta}}(\mathbf{s}, \mathbf{a})]$$

constant

???

$$Q^{\pi}(\mathbf{s}_t, \mathbf{a}_t) = E \left[\sum_{t'=t}^T \gamma^{t'-t} r(\mathbf{s}_{t'}, \mathbf{a}_{t'}) \right]$$

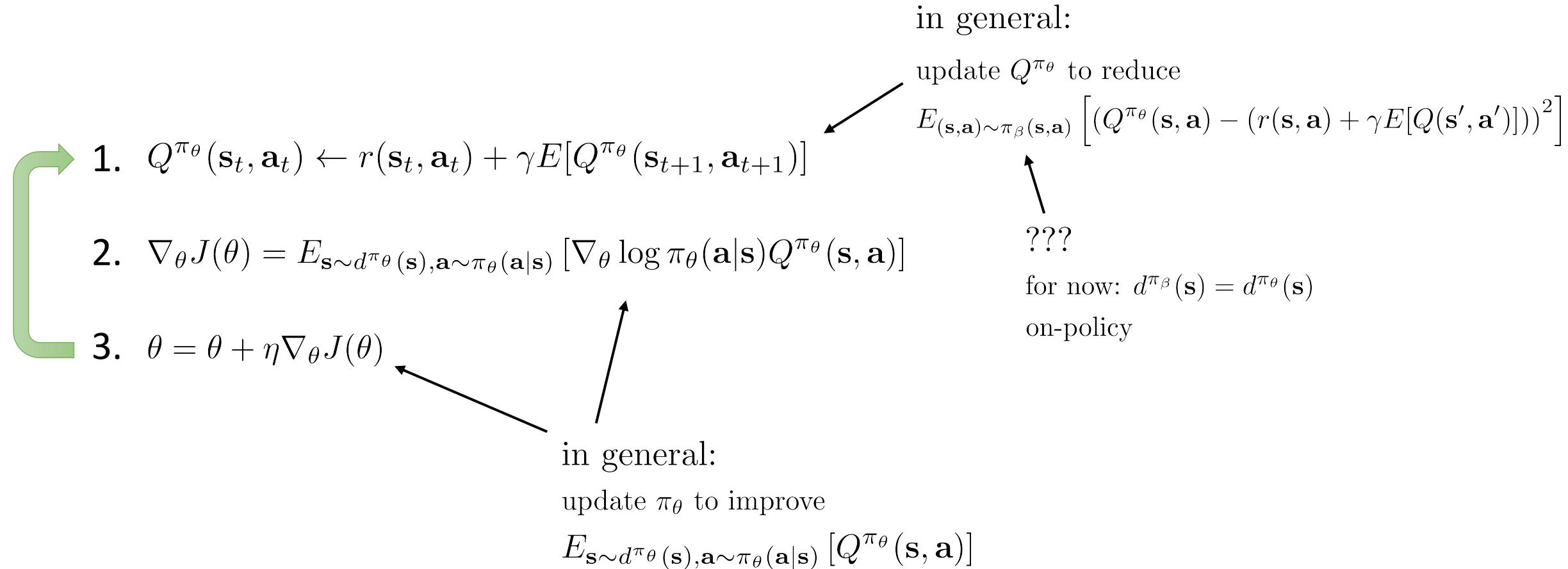
$$V^{\pi}(\mathbf{s}_t) = E \left[\sum_{t'=t}^T \gamma^{t'-t} r(\mathbf{s}_{t'}, \mathbf{a}_{t'}) \right]$$

these are really really useful!

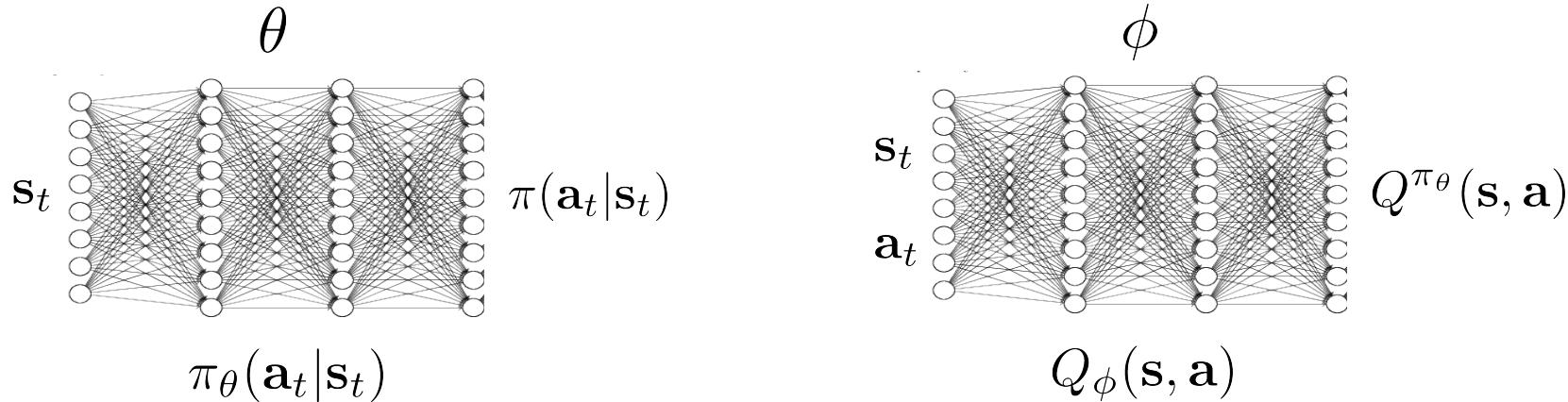
Bellman equations:

$$Q^{\pi}(\mathbf{s}_t, \mathbf{a}_t) = r(\mathbf{s}_t, \mathbf{a}_t) + \gamma E_{\mathbf{s}_{t+1} \sim p(\mathbf{s}_{t+1} | \mathbf{s}_t, \mathbf{a}_t), \mathbf{a}_{t+1} \sim \pi(\mathbf{a}_{t+1} | \mathbf{s}_{t+1})} [Q^{\pi}(\mathbf{s}_{t+1}, \mathbf{a}_{t+1})]$$

Basic actor-critic algorithm



General actor-critic with function approximation



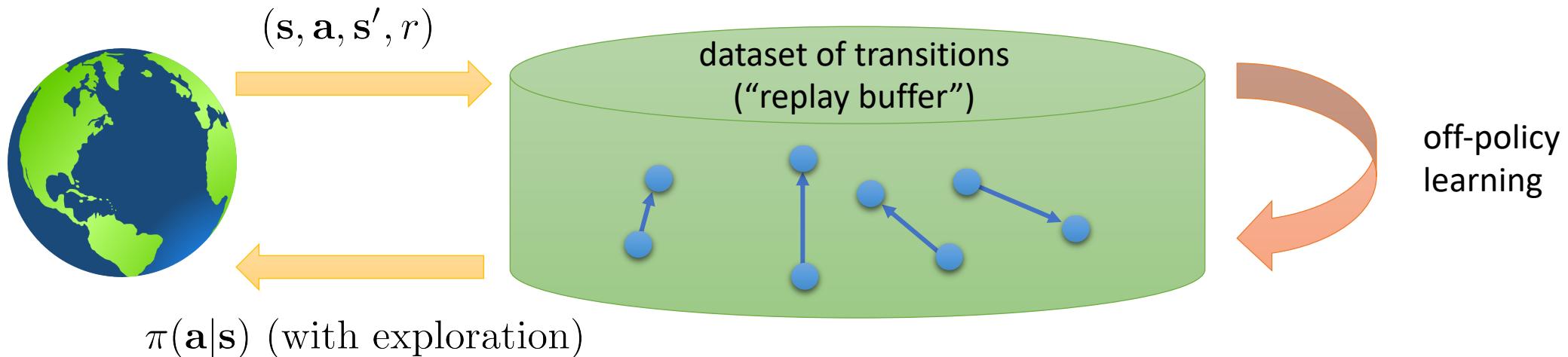
- 1. update Q_ϕ to decrease $E_{(s,a) \sim \pi_\beta(s,a)} \left[(Q_\phi(s, a) - (r(s, a) + \gamma E[Q_\phi(s', a')]))^2 \right]$
- 2. update π_θ to increase $E_{s \sim d^{\pi_\theta}(s), a \sim \pi_\theta(a|s)} [Q_\phi(s, a)]$

Switching to the *off-policy* setting

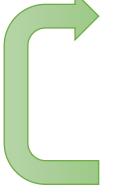
1. update Q_ϕ to decrease $E_{(\mathbf{s}, \mathbf{a}) \sim \pi_\beta(\mathbf{s}, \mathbf{a})} \left[(Q_\phi(\mathbf{s}, \mathbf{a}) - (r(\mathbf{s}, \mathbf{a}) + \gamma E[Q_\phi(\mathbf{s}', \mathbf{a}')]))^2 \right]$
2. update π_θ to increase $E_{\substack{\mathbf{s} \sim d^{\pi_\theta}(\mathbf{s}), \mathbf{a} \sim \pi_\theta(\mathbf{a}|\mathbf{s}) \\ d^{\pi_\beta}(\mathbf{s})}} [Q_\phi(\mathbf{s}, \mathbf{a})]$

is this still correct?

not exactly, but it works well in practice



Removing the actor (Q-learning/fitted Q-iteration)

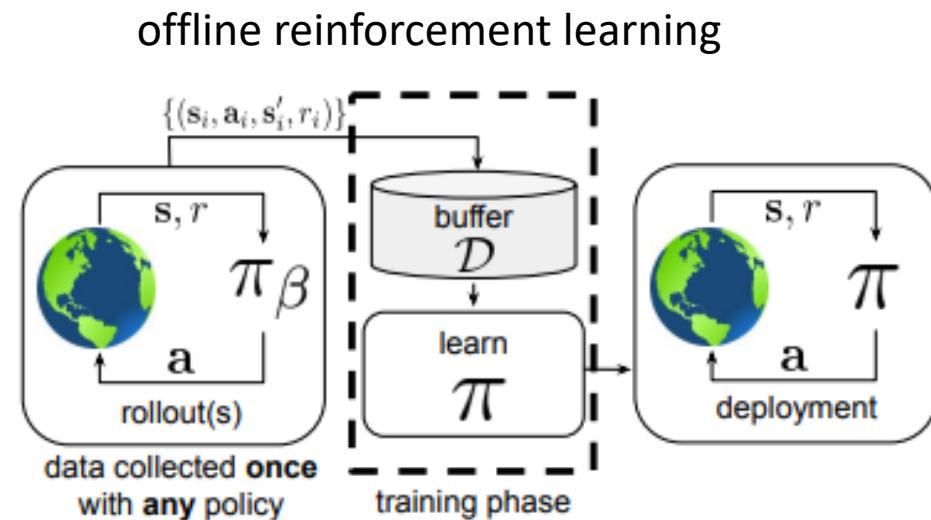
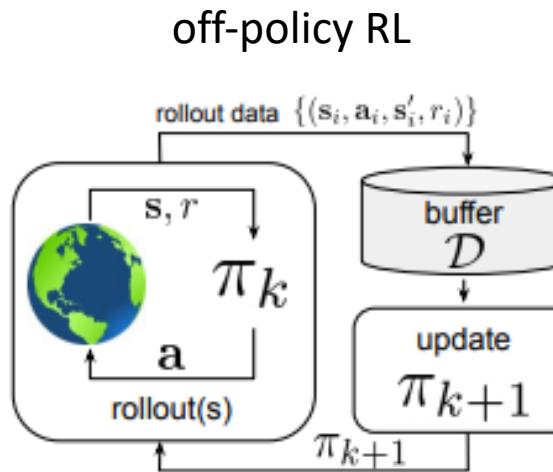
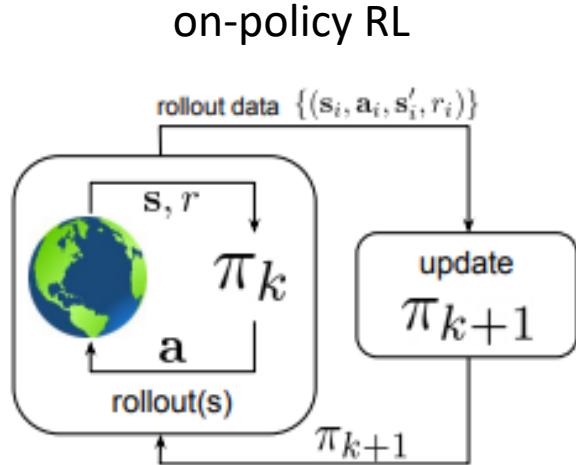
- 
1. update Q_ϕ to decrease $E_{(\mathbf{s}, \mathbf{a}) \sim \pi_\beta(\mathbf{s}, \mathbf{a})} \left[(Q_\phi(\mathbf{s}, \mathbf{a}) - (r(\mathbf{s}, \mathbf{a}) + \gamma E[Q_\phi(\mathbf{s}', \mathbf{a}')]))^2 \right]$
 2. ~~update π_θ to increase $E_{\mathbf{s} \sim d^{\pi_\beta}(\mathbf{s}), \mathbf{a} \sim \pi_\theta(\mathbf{a}|\mathbf{s})} [Q_\phi(\mathbf{s}, \mathbf{a})]$~~

choose π according to: $\pi(\mathbf{a}_t | \mathbf{s}_t) = \begin{cases} 1 & \text{if } \mathbf{a}_t = \arg \max_{\mathbf{a}_t} Q_\phi(\mathbf{s}_t, \mathbf{a}_t) \\ 0 & \text{otherwise} \end{cases}$

everything we discussed before still applies

1. Reinforcement learning primer
2. The offline reinforcement learning problem
3. Classic batch reinforcement learning algorithms
 - i. Importance-sampled policy gradients
 - ii. Least-squares temporal difference (LSTD)
4. Modern offline reinforcement learning algorithms
 - i. Policy constraint methods
 - ii. Model-based offline RL
 - iii. Value function regularization
 - iv. Uncertainty-based methods
 - v. Comparisons and discussion
5. Off-policy evaluation and model selection
 - i. Problem statement
 - ii. Direct methods and importance sampling methods
 - iii. Doubly robust methods
6. How should we evaluate offline RL methods?
 - i. When and why does offline RL work?
 - ii. How do we benchmark offline RL methods?
7. Conclusions, discussion, and open problems

What does offline RL mean?



Formally:

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

states are distributed by an unknown distribution d_{π}

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

$a \sim \pi_\beta(a|s)$ generally **not** known

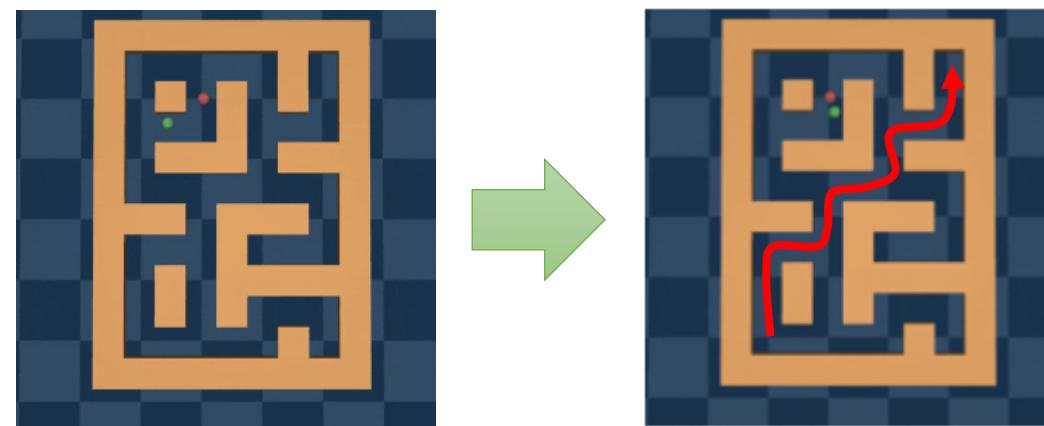
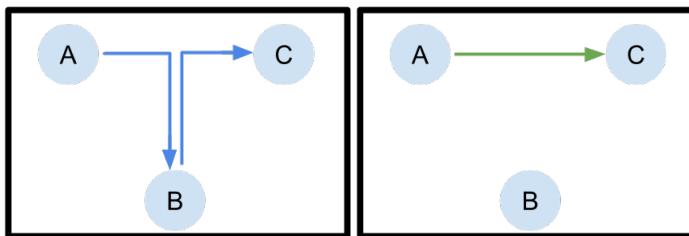
$s' \sim p(s'|s, a)$

$r \leftarrow r(s, a)$

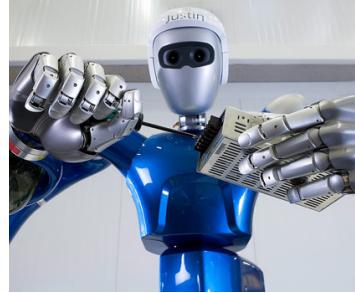
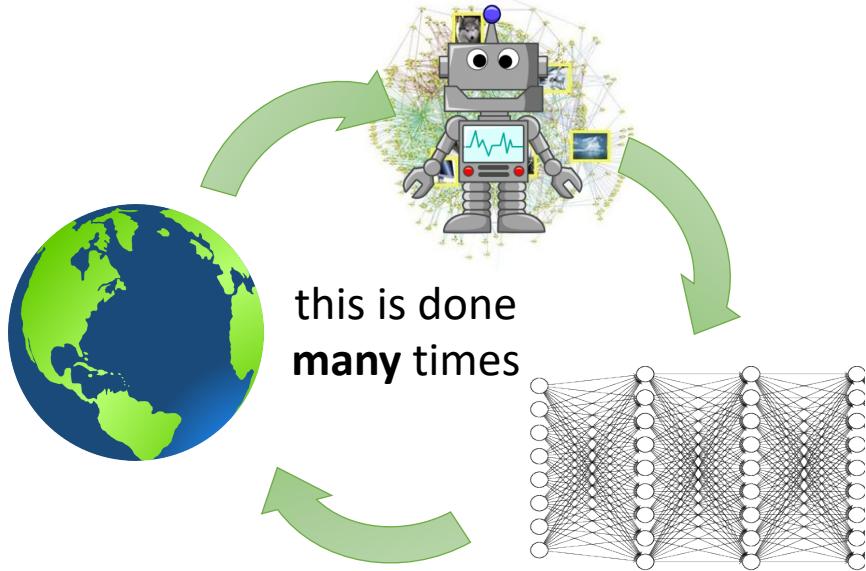
RL objective: $\max_{\pi} \sum_{t=0}^T E_{s_t \sim d^\pi(s), a_t \sim \pi(a|s)} [\gamma^t r(s_t, a_t)]$

How is this even possible?

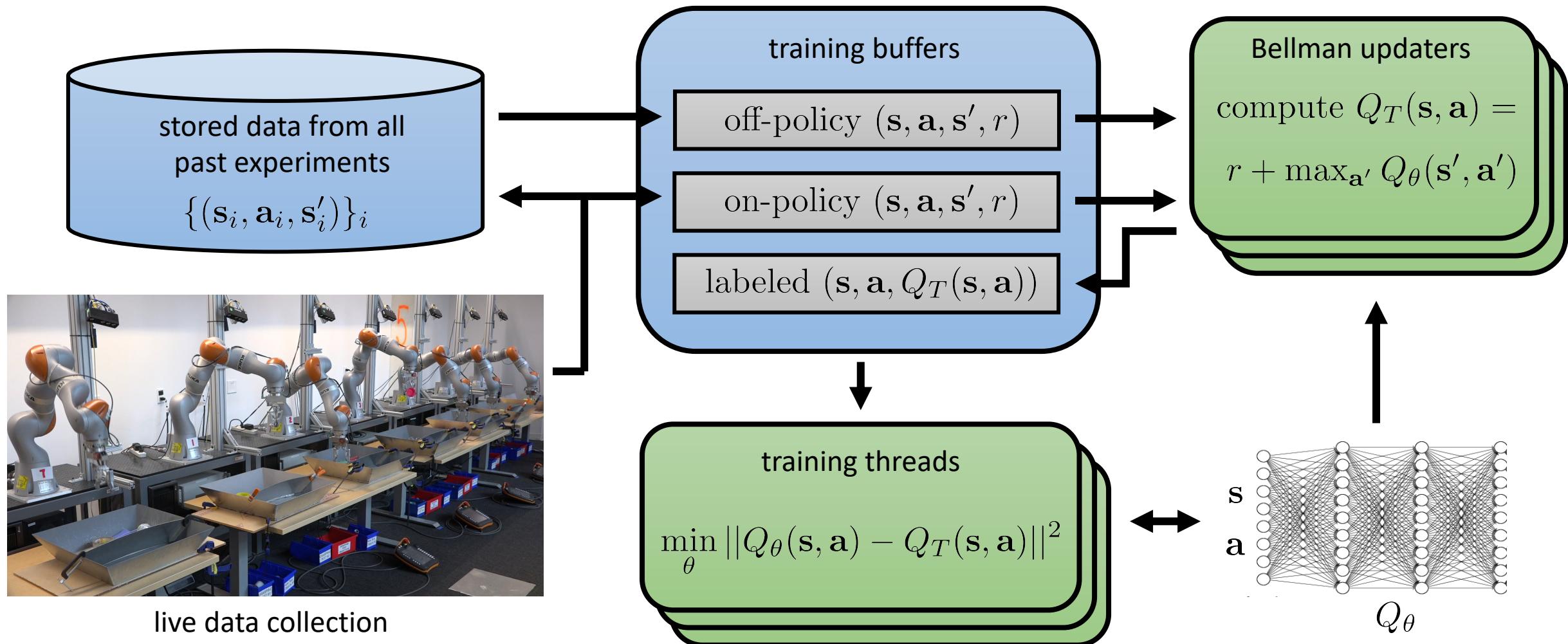
1. Find the “good stuff” in a dataset full of good and bad behaviors
2. Generalization: good behavior in one place may suggest good behavior in another place
3. “Stitching”: parts of good behaviors can be recombined



Why should we care?



Does it work?



Does it work?



Method	Dataset	Success	Failure
Offline QT-Opt	580k offline	87%	13%
Finetuned QT-Opt	580k offline + 28k online	96%	4%

Why is offline reinforcement learning hard?

Feedback: when we use a learned model/policy/etc. to act, it changes what we see

Simple example:

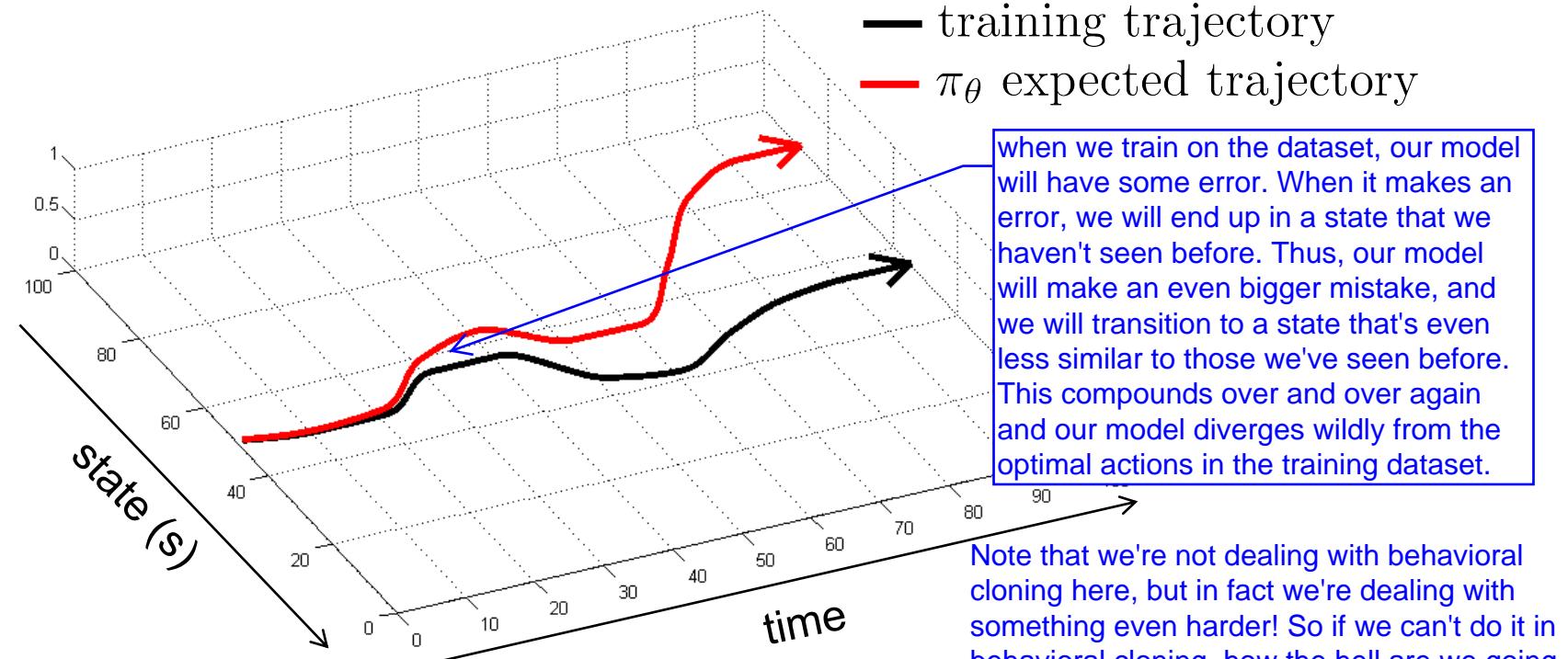
behavioral cloning

train π_θ to *copy*

assume data is *optimal*

assume we have a dataset with optimal actions. All we have to do is copy those optimal actions

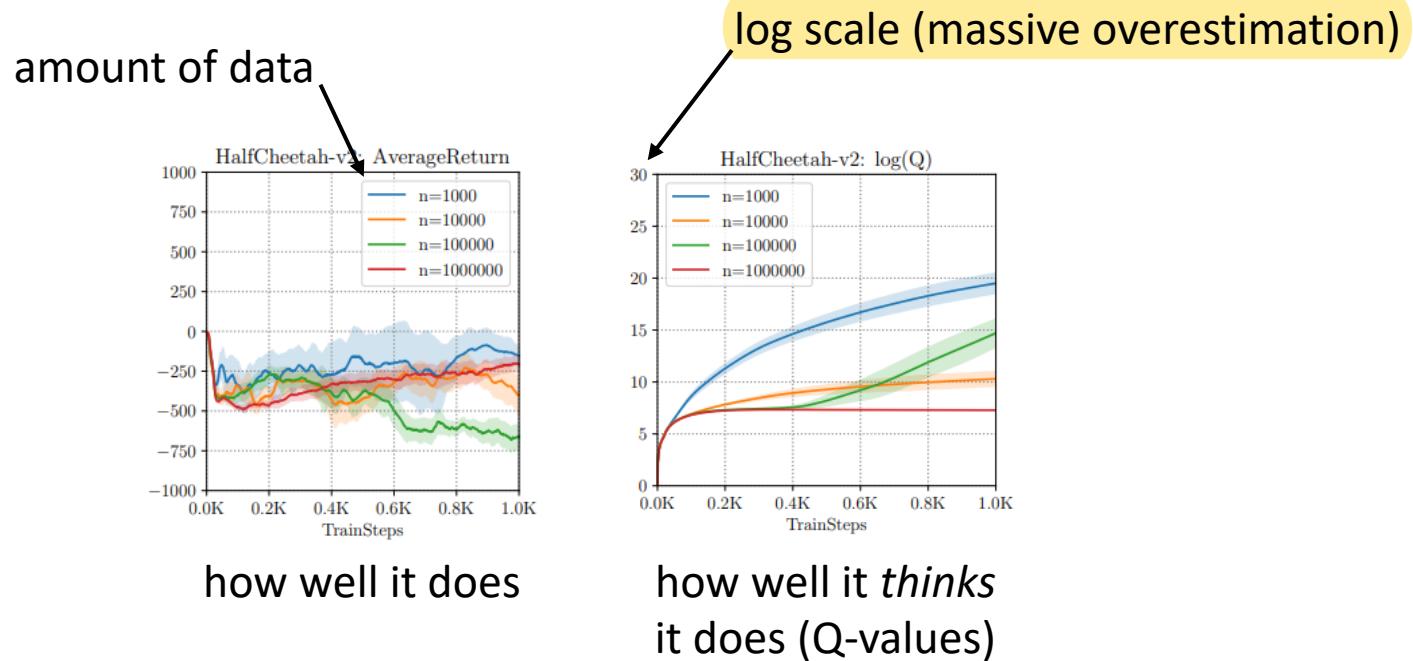
error scales as $O(T^2)$



Why is offline reinforcement learning hard?

Hypothesis 1: Overfitting

we overfit on the dataset we collect initially and bad stuff happens. They varied the amount of data available to the offline agent, and increasing the amount of data doesn't reduce the problem, so it seems that overfitting is not the issue.



Hypothesis 2: Training data is not good

Usually not the case: behavioral cloning of best data does better!

Distribution shift in a nutshell

Example empirical risk minimization (ERM) problem:

$$\theta \leftarrow \arg \min_{\theta} E_{\mathbf{x} \sim p(\mathbf{x}), y \sim p(y|\mathbf{x})} [(f_{\theta}(\mathbf{x}) - y)^2]$$

given some \mathbf{x}^* , is $f_{\theta}(\mathbf{x}^*)$ correct?

$$E_{\mathbf{x} \sim p(\mathbf{x}), y \sim p(y|\mathbf{x})} [(f_{\theta}(\mathbf{x}) - y)^2] \text{ is low}$$

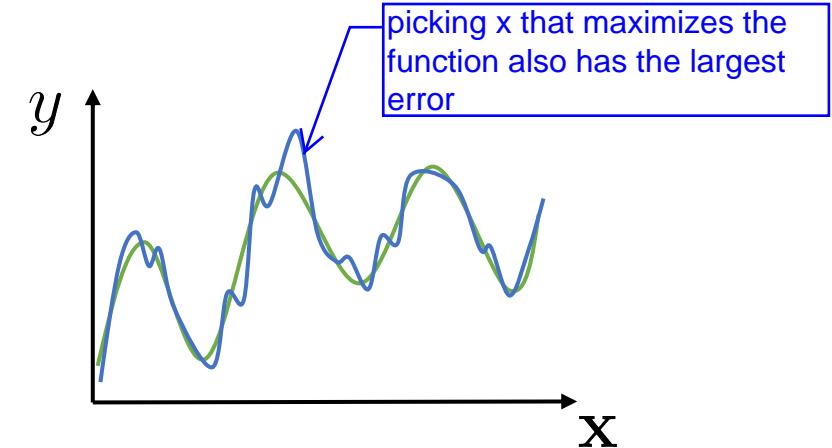
$$E_{\mathbf{x} \sim \bar{p}(\mathbf{x}), y \sim p(y|\mathbf{x})} [(f_{\theta}(\mathbf{x}) - y)^2] \text{ is not, for general } \bar{p}(\mathbf{x}) \neq p(\mathbf{x})$$

what if $\mathbf{x}^* \sim p(\mathbf{x})$? not necessarily...

usually we are not worried – neural nets generalize well!

what if we pick $\mathbf{x}^* \leftarrow \arg \max_{\mathbf{x}} f_{\theta}(\mathbf{x})$?

you can't expect the error to be low if \mathbf{x} was given by another distribution $\bar{p}(\mathbf{x})$



Where do we suffer from distribution shift?

The diagram illustrates the difference between the Q-learning setting and the Q-critic setting. In the top part, a blue box labeled "q-learning setting" has a blue arrow pointing down to the original Q-learning update rule:
$$Q(s, a) \leftarrow r(s, a) + \max_{a'} Q(s', a')$$
. A horizontal red line through this equation is crossed out with a red stroke. In the bottom part, a blue box labeled "Q-critic setting" has a blue arrow pointing up to the modified update rule:
$$Q(s, a) \leftarrow r(s, a) + E_{a' \sim \pi_{\text{new}}} [Q(s', a')]$$
. Below this, a bracket groups the term $E_{a' \sim \pi_{\text{new}}} [Q(s', a')]$ and the label $y(s, a)$, indicating they are equivalent in this setting.

what is the objective?

$$\min_Q E_{(\mathbf{s}, \mathbf{a}) \sim \pi_\beta(\mathbf{s}, \mathbf{a})} [(Q(\mathbf{s}, \mathbf{a}) - y(\mathbf{s}, \mathbf{a}))^2]$$



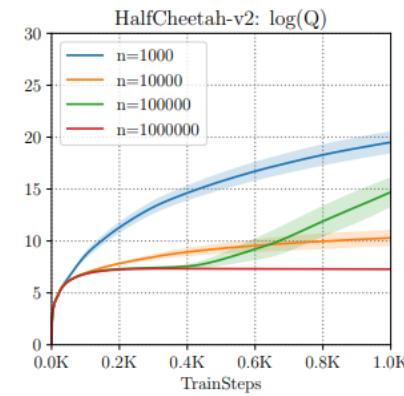
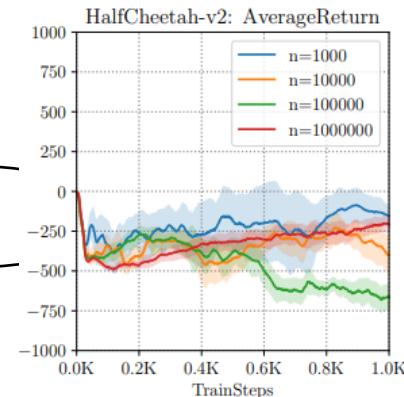
expect good accuracy when $\pi_\beta(\mathbf{a}|\mathbf{s}) = \pi_{\text{new}}(\mathbf{a}|\mathbf{s})$

how often does *that* happen?

even worse we pick π new that maximizes the Q-function

even worse: $\pi_{\text{new}} = \arg \max_{\pi} E_{\mathbf{a} \sim \pi(\mathbf{a}|\mathbf{s})}[Q(\mathbf{s}, \mathbf{a})]$

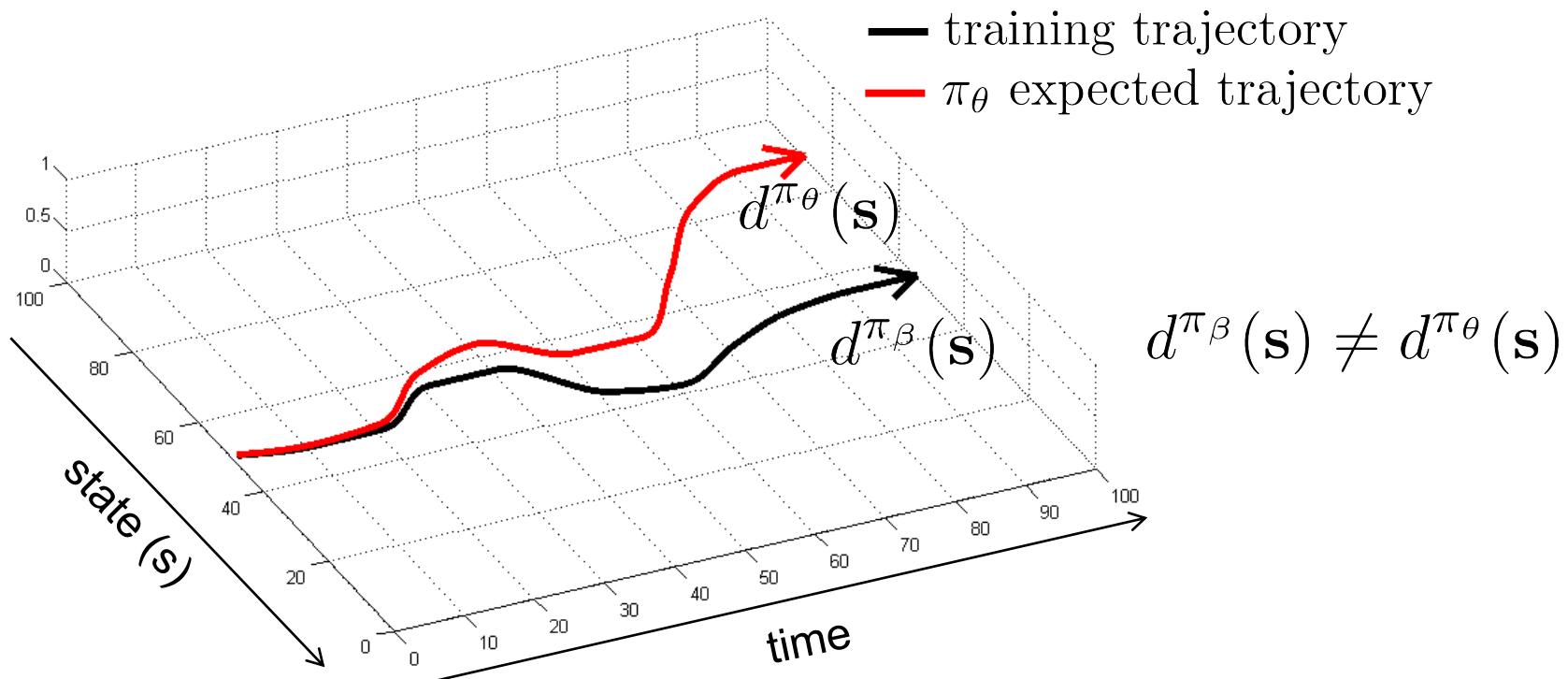
(what if we pick $\mathbf{x}^* \leftarrow \arg \max_{\mathbf{x}} f_{\theta}(\mathbf{x})$?).



how well it does

how well it *thinks*
it does (Q-values)

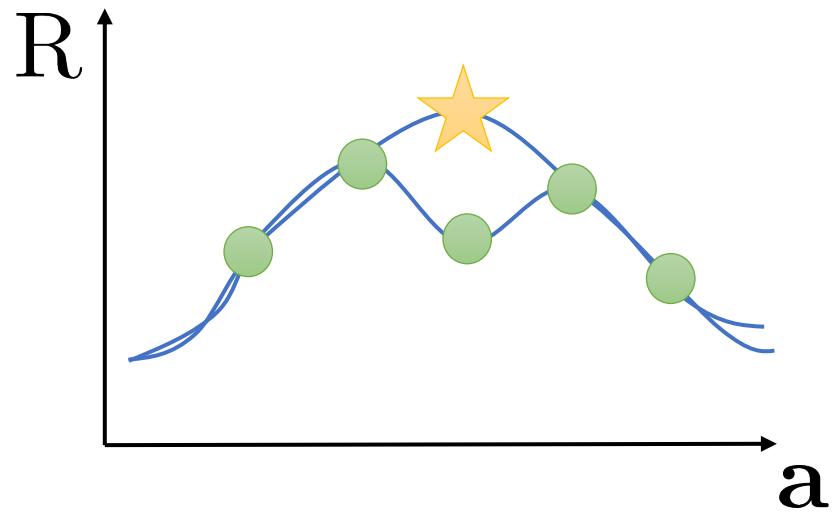
Where else do we suffer from distribution shift?



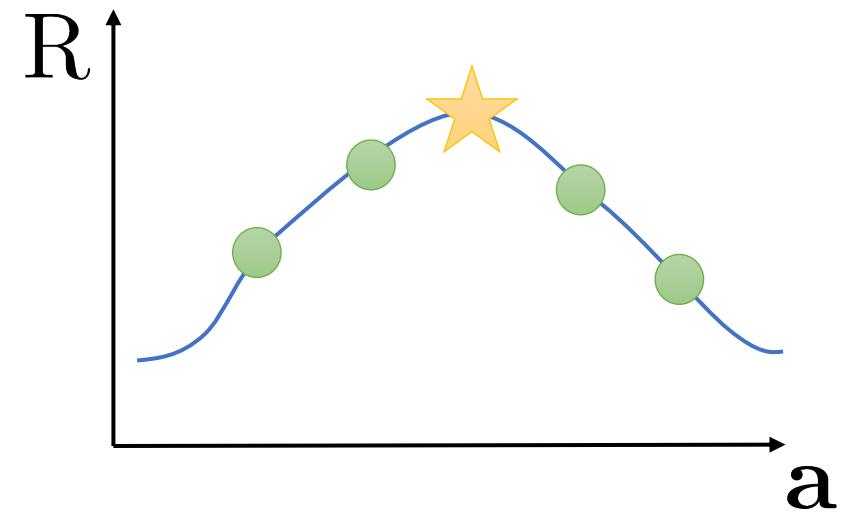
Even if π_θ is *great* on training data, it might be bad when used for control!

Sampling & function approximation error

online RL setting



offline RL setting



Existing challenges with sampling error and function approximation error in standard RL become **much more severe** in offline RL

1. Reinforcement learning primer
2. The offline reinforcement learning problem
3. **Classic batch reinforcement learning algorithms**
 - i. Importance-sampled policy gradients
 - ii. Least-squares temporal difference (LSTD)
4. Modern offline reinforcement learning algorithms
 - i. Policy constraint methods
 - ii. Model-based offline RL
 - iii. Value function regularization
 - iv. Uncertainty-based methods
 - v. Comparisons and discussion
5. Off-policy evaluation and model selection
 - i. Problem statement
 - ii. Direct methods and importance sampling methods
 - iii. Doubly robust methods
6. How should we evaluate offline RL methods?
 - i. When and why does offline RL work?
 - ii. How do we benchmark offline RL methods?
7. Conclusions, discussion, and open problems

Offline RL with policy gradients

RL objective: $\max_{\pi} \sum_{t=0}^T E_{\mathbf{s}_t \sim d^\pi(\mathbf{s}), \mathbf{a}_t \sim \pi(\mathbf{a}|\mathbf{s})} [\gamma^t r(\mathbf{s}_t, \mathbf{a}_t)]$

Trick when we want the expected value under one distribution but we have samples from another: Importance Sampling

$$\nabla_{\theta} J(\theta) = E_{\tau \sim \pi_{\theta}(\tau)} \left[\sum_{t=0}^T \nabla_{\theta} \gamma^t \log \pi_{\theta}(\mathbf{a}_t | \mathbf{s}_t) \hat{Q}(\mathbf{s}_t, \mathbf{a}_t) \right]$$

$$\approx \frac{1}{N} \sum_{i=1}^N \sum_{t=0}^T \nabla_{\theta} \gamma^t \log \pi_{\theta}(\mathbf{a}_{t,i} | \mathbf{s}_{t,i}) \hat{Q}(\mathbf{s}_{t,i}, \mathbf{a}_{t,i})$$

requires sampling from π_{θ} !

what if we only have samples from π_{β} ?

importance sampling:

$$\nabla_{\theta} J(\theta) \approx \frac{1}{N} \sum_{i=1}^N \underbrace{\frac{\pi_{\theta}(\tau_i)}{\pi_{\beta}(\tau_i)}}_{\text{importance weight}} \sum_{t=0}^T \nabla_{\theta} \gamma^t \log \pi_{\theta}(\mathbf{a}_{t,i} | \mathbf{s}_{t,i}) \hat{Q}(\mathbf{s}_{t,i}, \mathbf{a}_{t,i})$$

importance weight

high variance in practice

this is just a sum over future rewards

importance sampling

Offline RL with policy gradients

$$\nabla_{\theta} J(\theta) \approx \frac{1}{N} \sum_{i=1}^N \frac{\pi_{\theta}(\tau_i)}{\pi_{\beta}(\tau_i)} \sum_{t=0}^T \nabla_{\theta} \gamma^t \log \pi_{\theta}(\mathbf{a}_{t,i} | \mathbf{s}_{t,i}) \hat{Q}(\mathbf{s}_{t,i}, \mathbf{a}_{t,i})$$

$$\frac{\pi_{\theta}(\tau)}{\pi_{\beta}(\tau)} = \frac{p(s_1) \prod_t p(s_{t+1} | s_t, a_t) \pi_{\theta}(a_t | s_t)}{p(s_1) \prod_t p(s_{t+1} | s_t, a_t) \pi_{\beta}(a_t | s_t)}$$

can we fix this?

$$\nabla_{\theta} J(\theta) \approx \frac{1}{N} \sum_{i=1}^N \sum_{t=0}^T \underbrace{\left(\prod_{t'=0}^{t-1} \frac{\pi_{\theta}(a_{t',i} | s_{t',i})}{\pi_{\beta}(a_{t',i} | s_{t',i})} \right)}_{\text{accounts for difference in probability of landing in } s_{t,i}} \nabla_{\theta} \gamma^t \log \pi_{\theta}(\mathbf{a}_{t,i} | \mathbf{s}_{t,i}) \left(\prod_{t'=t}^T \frac{\pi_{\theta}(a_{t',i} | s_{t',i})}{\pi_{\beta}(a_{t',i} | s_{t',i})} \right) \hat{Q}(\mathbf{s}_{t,i}, \mathbf{a}_{t,i})$$

accounts for difference in probability of landing in $s_{t,i}$

we have $\mathbf{s}_t \sim d^{\pi_{\beta}}(\mathbf{s}_t)$, but want $\mathbf{s}_t \sim d^{\pi_{\theta}}(\mathbf{s}_t)$

why is this a reasonable approximation?
if π_{theta} lands you in similar states as π_{beta} , you can do alright by removing this. In math you can prove this

$$E_{\pi_{\theta}} \left[\sum_{t'=t}^T \gamma^{t'-t} r_{t'} \right] \approx \sum_{t'=t}^T \gamma^{t'-t} r_{t',i}$$

this is exponential in T
weights likely to be degenerate as T becomes large

we need to know
pi_beta

$$\left(\prod_{t'=t}^T \frac{\pi_{\theta}(a_{t',i} | s_{t',i})}{\pi_{\beta}(a_{t',i} | s_{t',i})} \right) \hat{Q}(\mathbf{s}_{t,i}, \mathbf{a}_{t,i})$$

accounts for having the incorrect $\hat{Q}(\mathbf{s}_{t,i}, \mathbf{a}_{t,i})$

corrects for the fact that \hat{Q} is an estimate of π_{β} not π_{θ}

Estimating the returns

recall that we want \hat{Q} to be an estimate of π_θ , not π_β . So we keep the importance sampling weights here.

$$\nabla_\theta J(\theta) \approx \frac{1}{N} \sum_{i=1}^N \sum_{t=0}^T \nabla_\theta \gamma^t \log \pi_\theta(\mathbf{a}_{t,i} | \mathbf{s}_{t,i}) \left(\prod_{t'=t}^T \frac{\pi_\theta(\mathbf{a}_{t',i} | \mathbf{s}_{t',i})}{\pi_\beta(\mathbf{a}_{t',i} | \mathbf{s}_{t',i})} \right) \hat{Q}(\mathbf{s}_{t,i}, \mathbf{a}_{t,i})$$



$$\sum_{t'=t}^T \left(\prod_{t''=t}^{t'} \frac{\pi_\theta(\mathbf{a}_{t'',i} | \mathbf{s}_{t'',i})}{\pi_\beta(\mathbf{a}_{t'',i} | \mathbf{s}_{t'',i})} \right) \gamma^{t'-t} r_{t',i}$$

but this is *still* exponential!

To avoid exponentially exploding importance weights, we **must** use value function estimation!

imagine we knew $Q^{\pi_\theta}(\mathbf{s}, \mathbf{a})$

We'll see how to do this shortly, but first let's conclude our discussion of importance sampling

The doubly robust estimator

$$\begin{aligned}
 V^{\pi_\theta}(\mathbf{s}_0) &\approx \sum_{t=0}^T \left(\prod_{t'=0}^{t'} \frac{\pi_\theta(\mathbf{a}_{tt'} | \mathbf{s}_{it} | \mathbf{s}_{t'}, i)}{\pi_\beta(\mathbf{a}_{tt'} | \mathbf{s}_{it} | \mathbf{s}_{t'}, i)} r_t \right) \gamma^{t'-t} r_{t',i} \\
 &= \sum_{t=0}^T \left(\prod_{t'=0}^t \rho_{t'} \right) \gamma^t r_t \\
 &= \rho_0 r_0 + \rho_0 \gamma \rho_1 r_1 + \rho_0 \gamma \rho_1 \gamma \rho_2 r_2 + \dots \\
 &= \rho_0 (r_0 + \gamma (\rho_1 (r_1 + \gamma (\rho_2 (r_2 + \gamma \dots))))) \\
 &= \bar{V}^T \quad \text{where } \bar{V}^{T+1-t} = \rho_t (r_t + \gamma \bar{V}^{T-t})
 \end{aligned}$$

this is exponential!

doubly robust estimation (bandit case)

$$V_{\text{DR}}(s) = \hat{V}(s) + \rho(s, a)(r_{s,a} - \hat{Q}(s, a))$$

model or function approximator

$$\bar{V}_{\text{DR}}^{T+1-t} = \hat{V}(\mathbf{s}_t) + \rho_t (r_t + \gamma \bar{V}_{\text{DR}}^{T-t} - \hat{Q}(\mathbf{s}_t, \mathbf{a}_t))$$

↓
model or function approximator

Marginalized importance sampling

Main idea: instead of using $\prod_t \frac{\pi_\theta(\mathbf{a}_t | \mathbf{s}_t)}{\pi_\beta(\mathbf{a}_t | \mathbf{s}_t)}$, estimate $w(\mathbf{s}, \mathbf{a}) = \frac{d^{\pi_\theta}(\mathbf{s}, \mathbf{a})}{d^{\pi_\beta}(\mathbf{s}, \mathbf{a})}$

if we can do this, we can estimate $J(\theta) \approx \frac{1}{N} \sum_i w(\mathbf{s}_i, \mathbf{a}_i) r_i$

typically this is done for off-policy evaluation, rather than policy learning

how to determine $w(\mathbf{s}, \mathbf{a})$? typically solve some kind of consistency condition

example (Zhang et al., GenDICE):

$$d^{\pi_\beta}(\mathbf{s}', \mathbf{a}') w(\mathbf{s}', \mathbf{a}') = \underbrace{(1-\gamma)p_0(\mathbf{s}')\pi_\theta(\mathbf{a}'|\mathbf{s}')}_{\text{probability of starting in } (\mathbf{s}', \mathbf{a}')} + \gamma \underbrace{\sum_{\mathbf{s}, \mathbf{a}} \pi_\theta(\mathbf{a}'|\mathbf{s}') p(\mathbf{s}'|\mathbf{s}, \mathbf{a}) d^{\pi_\beta}(\mathbf{s}, \mathbf{a}) w(\mathbf{s}, \mathbf{a})}_{\text{probability of transitioning into } (\mathbf{s}', \mathbf{a}')}$$

solving for $w(\mathbf{s}, \mathbf{a})$ typically involves some fixed point problem

Additional readings: importance sampling

Classic work on importance sampled policy gradients and return estimation:

Precup, D. (2000). Eligibility traces for off-policy policy evaluation.

Peshkin, L. and Shelton, C. R. (2002). Learning from scarce experience.

Doubly robust estimators and other improved importance-sampling estimators:

Jiang, N. and Li, L. (2015). Doubly robust off-policy value evaluation for reinforcement learning.

Thomas, P. and Brunskill, E. (2016). Data-efficient off-policy policy evaluation for reinforcement learning.

Analysis and theory:

Thomas, P. S., Theocharous, G., and Ghavamzadeh, M. (2015). High-confidence off-policy evaluation.

Marginalized importance sampling:

Hallak, A. and Mannor, S. (2017). Consistent on-line off-policy evaluation.

Liu, Y., Swaminathan, A., Agarwal, A., and Brunskill, E. (2019). Off-policy policy gradient with state distribution correction.

Additional readings in our offline RL survey: Section 3.1, 3.2, 3.3, 3.4: <https://arxiv.org/abs/2005.01643>

Offline value function estimation

How have people thought about it before?

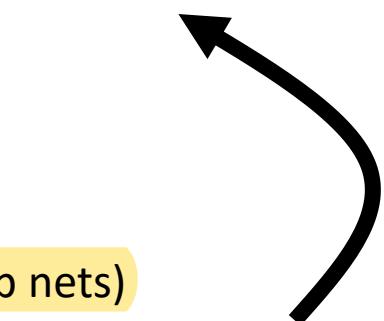
Extend existing ideas for approximate dynamic programming and Q-learning to offline setting

Derive tractable solutions with simple (e.g., linear) function approximators

How are people thinking about it now?

Derive approximate solutions with highly expressive function approximators (e.g., deep nets)

The primary challenge turns out to be **distributional shift**



generally not concerned with distributional shift before

(maybe it was not such a big problem with linear models)

We'll discuss some older offline/batch RL methods
next for completeness

Warmup: linear models

Φ – feature matrix, $|S| \times K$

could also think of as a vector-valued function $\Phi(\mathbf{s})$

Can we do (offline) model-based RL in feature space?

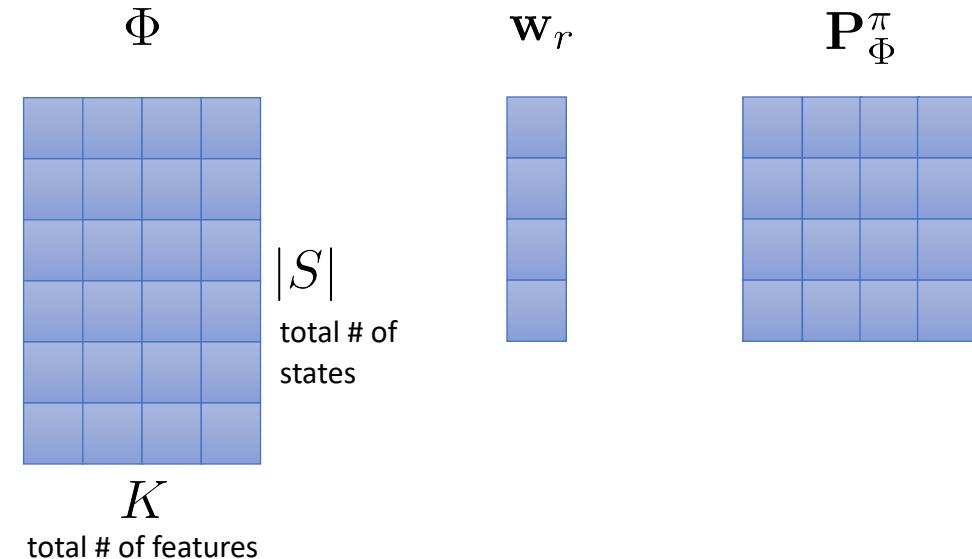
1. Estimate the reward
2. Estimate the transitions
3. Recover the value function
4. Improve the policy

$$1. \text{ Reward model: } \Phi \mathbf{w}_r \approx r$$

$$2. \text{ Transition model: } \Phi \mathbf{P}_\Phi \approx \mathbf{P}^\pi \Phi$$

estimated feature-space
transition matrix
 $K \times K$

real transition matrix
(on states)
 $|S| \times |S|$



vector of rewards for all state-action tuples
but we'll talk about sample-based setting soon!

$$\text{least squares: } \mathbf{w}_r = (\Phi^T \Phi)^{-1} \Phi^T \vec{r}$$

$$\text{least squares: } \mathbf{P}_\Phi = (\Phi^T \Phi)^{-1} \Phi^T \mathbf{P}^\pi \Phi$$

all of this is for a *fixed* policy π

Recovering the value function

1. Reward model: $\Phi \mathbf{w}_r \approx r$

$$\text{least squares: } \mathbf{w}_r = (\Phi^T \Phi)^{-1} \Phi \vec{r}$$

2. Transition model: $\Phi \mathbf{P}_\Phi \approx \mathbf{P}^\pi \Phi$

$$\text{least squares: } \mathbf{P}_\Phi = (\Phi^T \Phi)^{-1} \Phi \mathbf{P}^\pi \Phi$$

3. Estimate $V^\pi \approx V_\Phi^\pi = \Phi \mathbf{w}_V$

can apply the same equation in feature space:

$$\mathbf{w}_V = (\mathbf{I} - \gamma \mathbf{P}_\Phi)^{-1} \mathbf{w}_r$$

← substitute

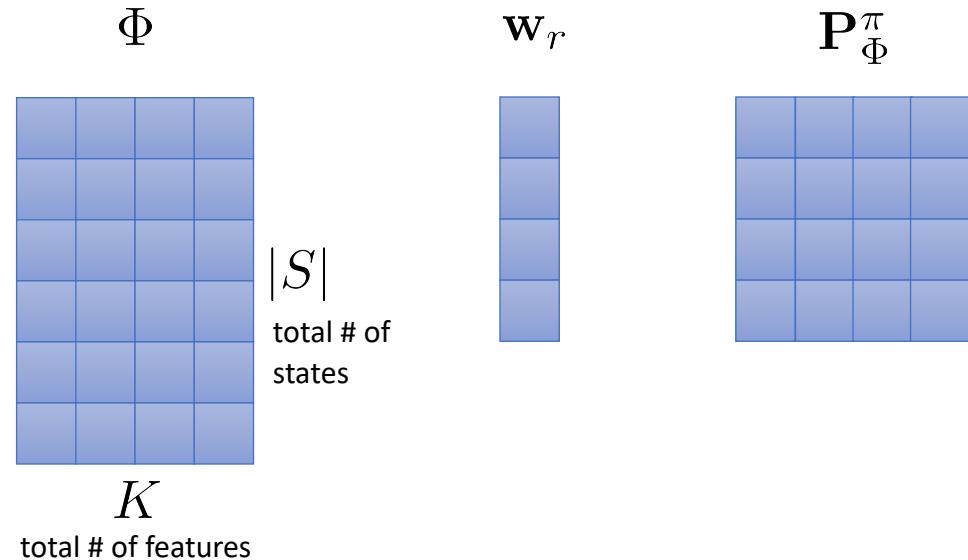
but wait – do we even *need* the model?

$$\mathbf{w}_V = (\mathbf{I} - \gamma(\Phi^T \Phi)^{-1} \Phi^T \mathbf{P}^\pi \Phi)^{-1} (\Phi^T \Phi)^{-1} \Phi^T \vec{r}$$

after a bit of algebra...

$$\mathbf{w}_V = (\Phi^T \Phi - \gamma \Phi^T \mathbf{P}^\pi \Phi)^{-1} \Phi^T \vec{r}$$

this is called least-squares temporal difference (LSTD)



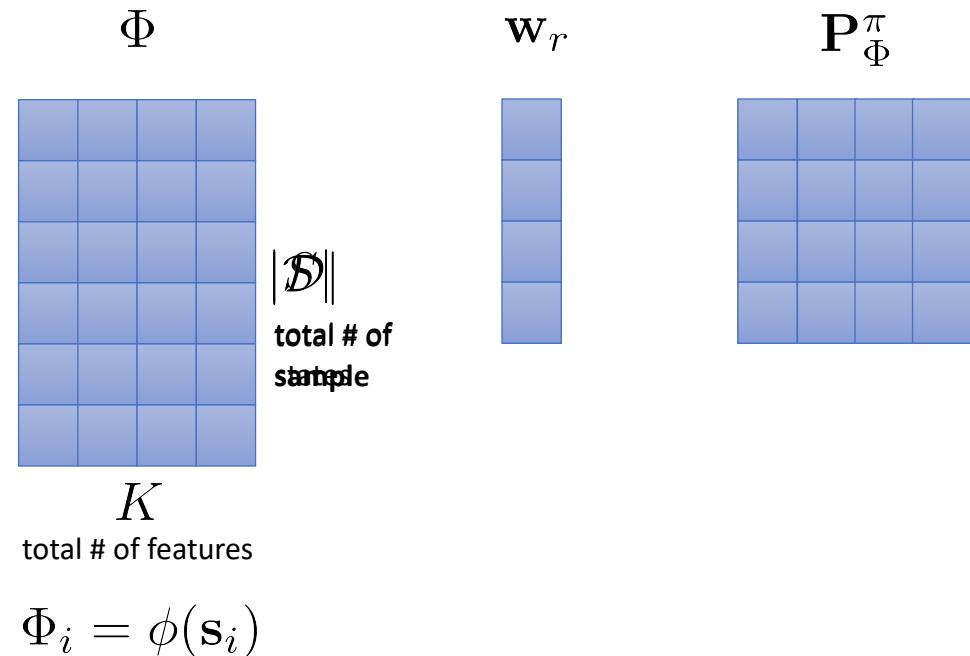
Aside: solving for V^π in terms of \mathbf{P}^π and \mathbf{r} :

$$\begin{aligned} V^\pi &= \mathbf{r} + \gamma \mathbf{P}^\pi V^\pi \\ (\mathbf{I} - \gamma \mathbf{P}^\pi) V^\pi &= \mathbf{r} \\ V^\pi &= (\mathbf{I} - \gamma \mathbf{P}^\pi)^{-1} \mathbf{r} \end{aligned}$$

Doing it all with samples

$$\mathbf{w}_V = (\Phi^T \Phi - \gamma \Phi^T \mathbf{P}^\pi \Phi)^{-1} \Phi^T \vec{\mathbf{r}}$$

$\mathcal{D} = \{(\mathbf{s}_i, \mathbf{a}_i, r_i, \mathbf{s}'_i)\}$ ↑ $\vec{\mathbf{r}}_i = r(\mathbf{s}_i, \mathbf{a}_i)$
 replace with Φ'
 $\Phi'_i = \phi(\mathbf{s}'_i)$



Everything else works **exactly** the same way, only now we have some sampling error

Improving the policy

1. Estimate the reward
2. Estimate the transitions
3. Recover the value function
4. Improve the policy

or just do these together with LSTD!

typical policy improvement step:

~~→ $\pi'(s) \leftarrow \text{Greedy}(\Phi w_V)$~~

~~estimate $V^{\pi'}$~~

That's not going to work for offline RL!

$$w_V = (\Phi^T \Phi - \gamma \Phi^T \cancel{P^\pi} \Phi)^{-1} \Phi^T \vec{r}$$
$$\mathcal{D} = \{(s_i, a_i, r_i, s'_i)\}$$

↑ ↑
 $\vec{r}_i = r(s_i, a_i)$

replace with Φ'
 $\Phi'_i = \phi(s'_i)$

this requires samples from π !

Least-squares policy iteration (LSPI)

Main idea: replace LSTD with LSTDQ – LSTD but for Q-functions

$$\mathbf{w}_Q = (\Phi^T \Phi - \gamma \Phi^T \Phi')^{-1} \Phi^T \vec{\mathbf{r}}$$

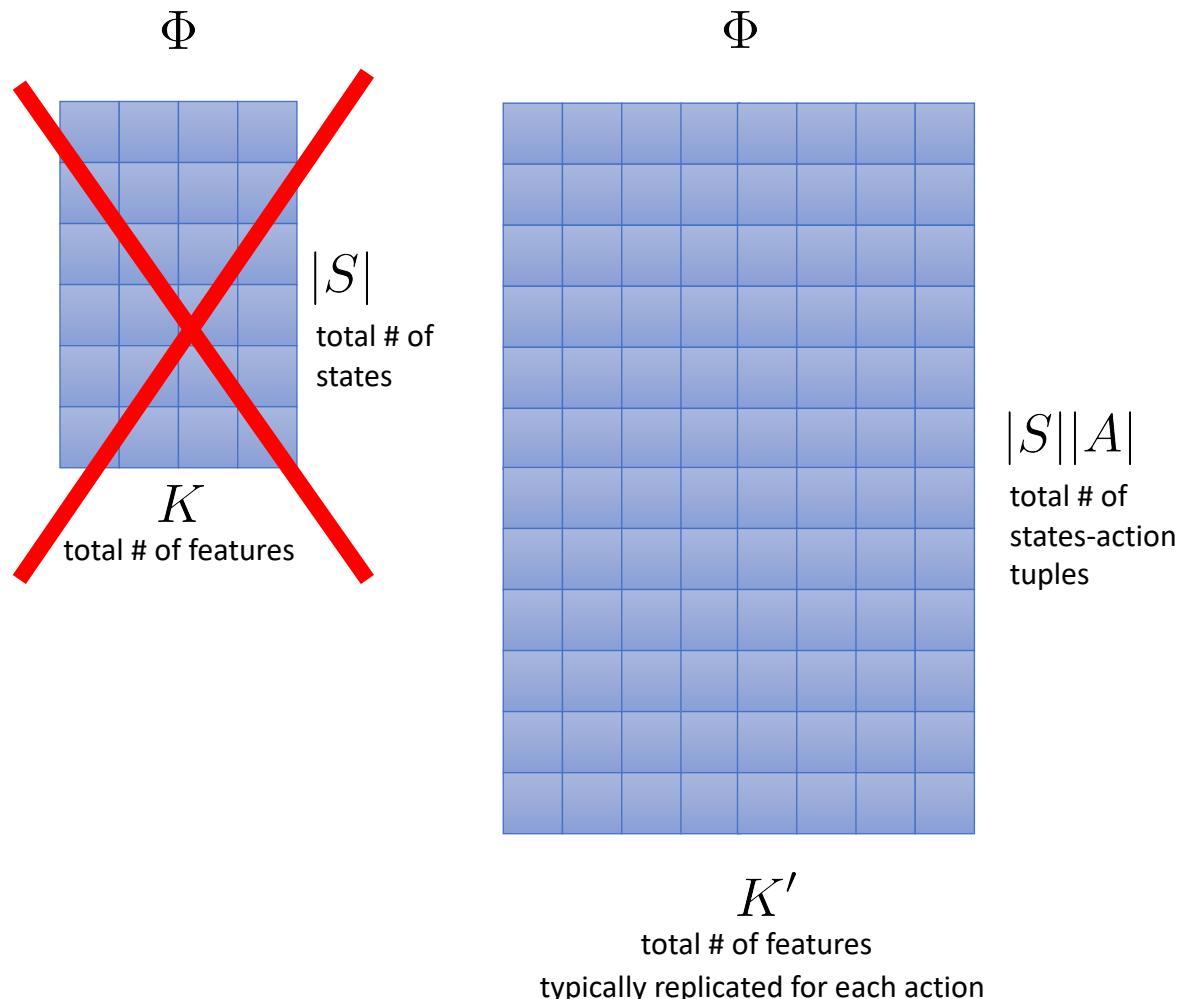
$$\mathcal{D} = \{(\mathbf{s}_i, \mathbf{a}_i, r_i, \mathbf{s}'_i)\}$$

$$\Phi'_i = \phi(\mathbf{s}'_i, \pi(\mathbf{s}'_i))$$

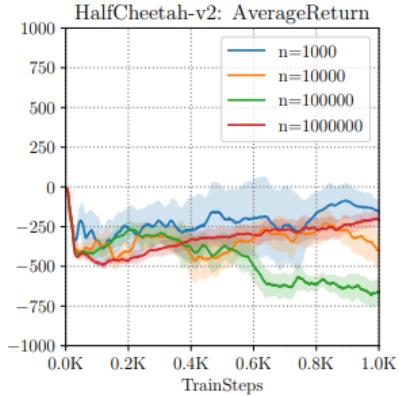
LSPI:

1. compute \mathbf{w}_Q for π_k
2. $\pi_{k+1}(\mathbf{s}) = \arg \max_{\mathbf{a}} \phi(\mathbf{s}, \mathbf{a}) \mathbf{w}_Q$
3. Set $\Phi'_i = \phi(\mathbf{s}'_i, \pi_{k+1}(\mathbf{s}'_i))$

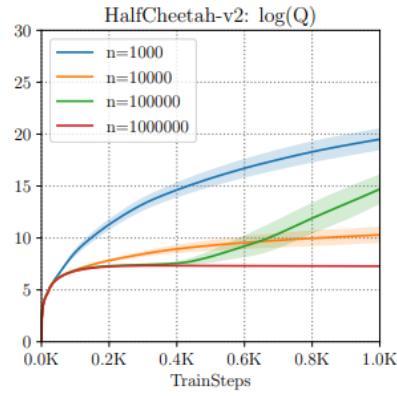
encode the action π would take
not the action in the data



What's the issue?



how well it does



how well it *thinks* it does (Q-values)

$$Q(s, a) \leftarrow \underbrace{r(s, a) + E_{a' \sim \pi_{\text{new}}}[Q(s', a')]}_{y(s, a)}$$

expect good accuracy when $\pi_\beta(a|s) = \pi_{\text{new}}(a|s)$

even *worse*: $\pi_{\text{new}} = \arg \max_\pi E_{a \sim \pi(a|s)}[Q(s, a)]$

In general, all approximate dynamic programming (e.g., fitted value/Q iteration) methods will suffer from action distributional shift, and we **must** fix it!

$$\min_Q E_{(s,a) \sim \pi_\beta(s,a)} [(Q(s, a) - y(s, a))^2]$$

↑
behavior policy ↑
target value

how often does *that* happen?

1. Reinforcement learning primer
2. The offline reinforcement learning problem
3. Classic batch reinforcement learning algorithms
 - i. Importance-sampled policy gradients
 - ii. Least-squares temporal difference (LSTD)
- 4. Modern offline reinforcement learning algorithms**
 - i. Policy constraint methods
 - ii. Model-based offline RL
 - iii. Value function regularization
 - iv. Uncertainty-based methods
 - v. Comparisons and discussion
5. Off-policy evaluation and model selection
 - i. Problem statement
 - ii. Direct methods and importance sampling methods
 - iii. Doubly robust methods
6. How should we evaluate offline RL methods?
 - i. When and why does offline RL work?
 - ii. How do we benchmark offline RL methods?
7. Conclusions, discussion, and open problems

Addressing Action Distributional Shift in Offline RL

$$Q(s, a) \leftarrow r(s, a) + \gamma \mathbb{E}_{a' \sim \pi_\theta(a'|s')} [Q(s', a')]$$

so we're effectively sampling alphas from the π_θ distribution

a' can be OOD from π_β

so when we do the bellman backup, we backup values that are erroneous and not in the dataset, because π_β was used to generate the data!

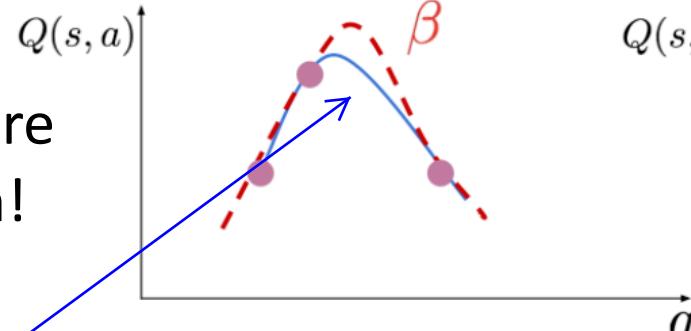
$$\pi_\theta := \arg \max_{\pi_\theta} \mathbb{E}_{s \sim \mathcal{D}, a \sim \pi_\theta(a|s)} [Q(s, a)] \text{ s.t. } D(\pi_\theta(a|s), \pi_\beta(a|s)) \leq \varepsilon$$

"Policy constraint"

Limits OOD actions used for training the Q-function.

Actions used for training are also (almost) trained upon!

red line represents behavior policy

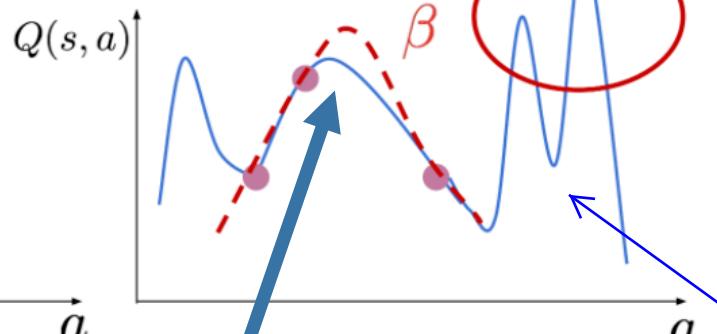


blue represents the Q-function

when we use a policy constraint, we ensure that the chosen action a' is close to the distribution of the behavior policy. This prevents the action distribution shift

$$\mathbb{E}_{a' \sim \pi_\theta(a'|s')} [Q(s', a')]$$

Will be used for backups



these actions are OOD for π_β . As such, our Q-value is bad at estimating them. Some of our estimates will be erroneously high. But using the bellman backup, we will pick the actions that maximize these!

Different Types of Policy Constraint Methods

$$\pi_\theta := \arg \max_{\pi_\theta} \mathbb{E}_{s \sim \mathcal{D}, a \sim \pi_\theta(a|s)} [Q(s, a)] \text{ s.t. } D(\pi_\theta(a|s), \pi_\beta(a|s)) \leq \varepsilon$$

Distribution constraints

(Fujimoto et al. 2019: action clipping)
 (Wu et al. 2019, Jaques et al. 2019)

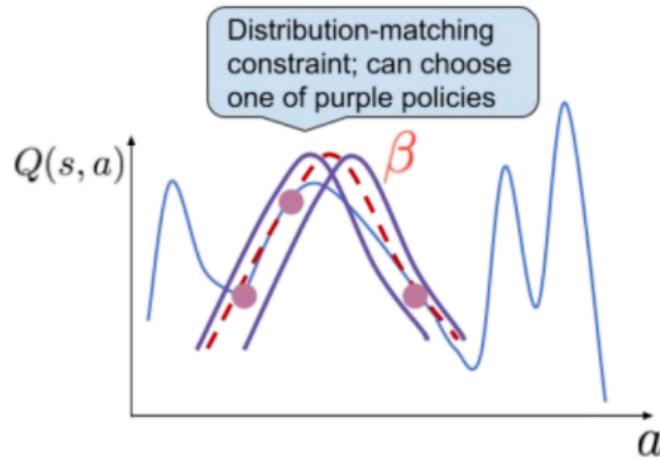
$$D(\pi_\theta, \pi_\beta) = D_{\text{KL}}(\pi_\theta, \pi_\beta)$$

KL divergence

$$D(\pi_\theta, \pi_\beta) = D_f(\pi_\theta, \pi_\beta) = \mathbb{E}_{\pi_\beta} \left[f \left(\frac{\pi_\theta}{\pi_\beta} \right) \right]$$

f-divergence

(Wu et al. 2019)



(Peng et al. 2019, Seigel et al. 2019,
 Wang et al. 2020, Nair et al. 2020)

Several other forms. Different benefits over each other

$$\max_{\pi_\theta(a|s)} \mathbb{E}_{\pi_\theta} [Q(s, a)] - \alpha D_{\text{KL}}(\pi_\theta(a|s), \pi_\beta(a|s))$$



(solve in closed form)

$$\pi_\theta(a|s) \propto \pi_\beta(a|s) \exp(Q(s, a)/\alpha)$$

you can implement this in code

“Implicit” distributional constraint

Different Types of Policy Constraint Methods

$$\pi_\theta := \arg \max_{\pi_\theta} \mathbb{E}_{s \sim \mathcal{D}, a \sim \pi_\theta(a|s)} [Q(s, a)] \quad \text{s.t. } D(\pi_\theta(a|s), \pi_\beta(a|s)) \leq \varepsilon$$

Support constraints

$$D(\pi_\theta, \pi_\beta) = \mathbf{D}(\pi_\theta, \text{supp}(\pi_\beta))$$

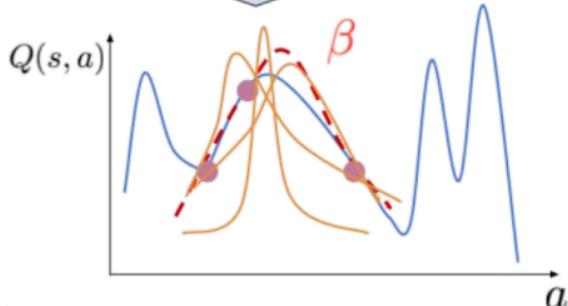
constrain the learned policy π_θ to the support of π_β

Can also use other integral probability metrics

(Sriperumbudur et al. 2009)

$$\text{IPM}_\mathcal{F}(\pi_\theta, \pi_\beta) = \max_{f \in \mathcal{F}} |\mathbb{E}_{\pi_\theta}[f(s, a)] - \mathbb{E}_{\pi_\beta}[f(s, a)]|$$

Support constraint;
can choose one of
yellow policies



Discrete Action Spaces

$$D(\pi_\theta, \pi_\beta) = \sum_{a \notin \mathcal{D}[s]} \pi_\theta(a|s)$$

Support at state s

$$Q(s, a) \leftarrow r(s, a) + \gamma \max_{a'} Q(s', a')$$

$$D(\pi_\theta, \pi_\beta) = \text{MMD}(\pi_\theta, \pi_\beta)$$

(Kumar et al. 2019)

Maximum mean discrepancy

(Gretton et al. 2012)

we can implement this constraint by calculating the MMD

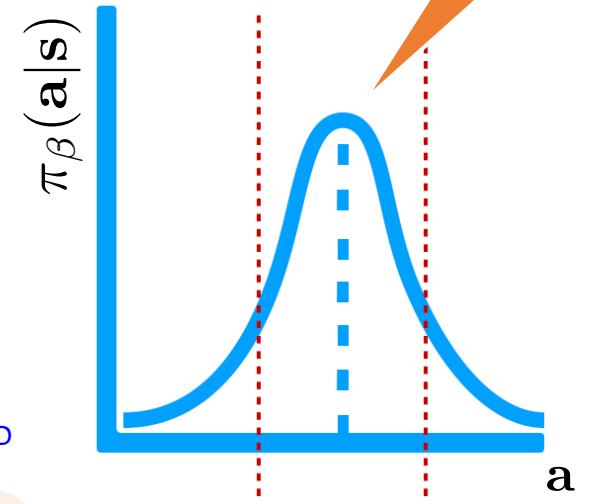
$$\text{Constrain } \pi_\theta(a|s) = \pi_\beta(a|s)$$

on out-of-support actions
(Laroche et al. 2019)

$$\max_{a'} Q(s', a')$$

$$\max_{a' \in \mathcal{D}[s']} Q(s', a')$$

Discrete BCQ (Fujimoto et al. 2019),
EMaQ (Ghasemipour et al. 2020)



Intuition: compare samples
regardless of density

we don't really care
about the densities

the "support" of the behavior policy
is the set of all possible actions that
are highly likely under the behavior
policy

restrict the support of the actions
that you compute the max over!

Comparison of Policy Constraint Methods

Do different policy constraint methods behave differently?

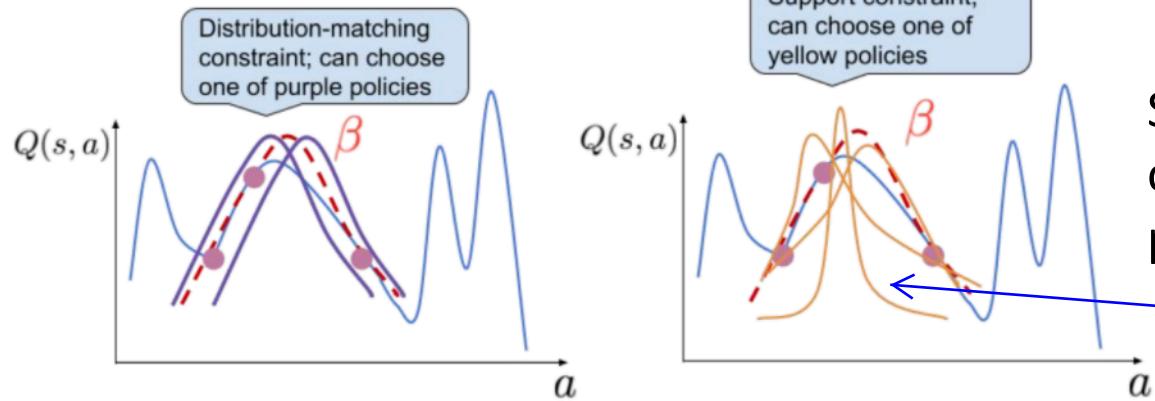
$$\max_{\pi} \mathbb{E}_{s_t \sim d^\pi(s), a_t \sim \pi(a|s)} \left[\sum_t \gamma^t r(s_t, a_t) \right] - \alpha D(\pi(a|s), \pi_\beta(a|s))$$

Standard RL objective Conservative behavior

Different constraints alter the behavior of the algorithm differently!

Support vs. distribution constraints

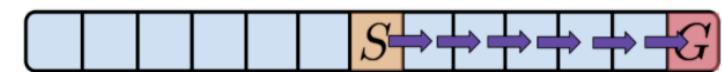
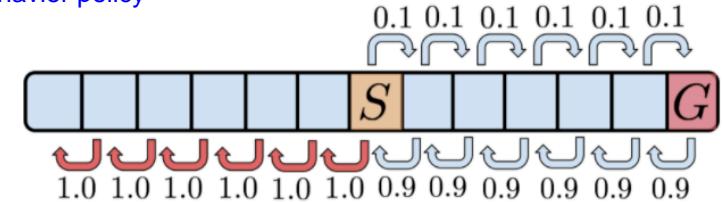
Better in theory, since more flexibility to choose policies



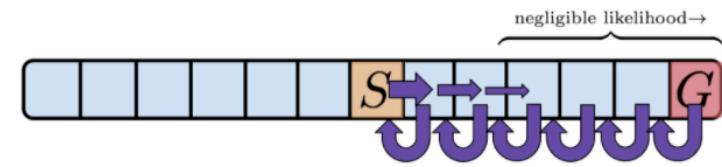
Support constraints still allow convergence to the optimal policy with full coverage

notice that the densities are really different, but they all are within the support of π_β

policy constraint methods alter the optimization function. We're no longer solving the traditional RL objective because we introduce another term penalizing us for deviating too far from the behavior policy



(b) Learned Policy via support-constraint



(a) Learned Policy via distribution-matching

so far we've looked at constraining π_{β} and π_{θ} . So we look at $\pi_{\theta}(a|s)$ and compare it to $\pi_{\beta}(a|s)$

A different variant of policy constraints

So far, we have discussed constraints on the policy distribution.

Can also define constraints on state-action marginal distributions

$$D(\pi_{\theta}, \pi_{\beta}) = D_{\text{KL}}(d^{\pi}(s, a), d^{\pi_{\beta}}(s, a))$$

Also a “policy constraint” since we technically constrain policies

marginal state-action distribution under π_{β}

Primarily explored for off-policy evaluation, but can also be used for offline RL

Lots of work:

“DICE” methods: DualDICE (Nachum et al. 2019), GenDICE (Zhang et al. 2020), AlgaeDICE (Nachum et al. 2019), GradientDICE (Zhang et al. 2020) etc.

Minimax weight learning: Uehara and Jiang, 2019

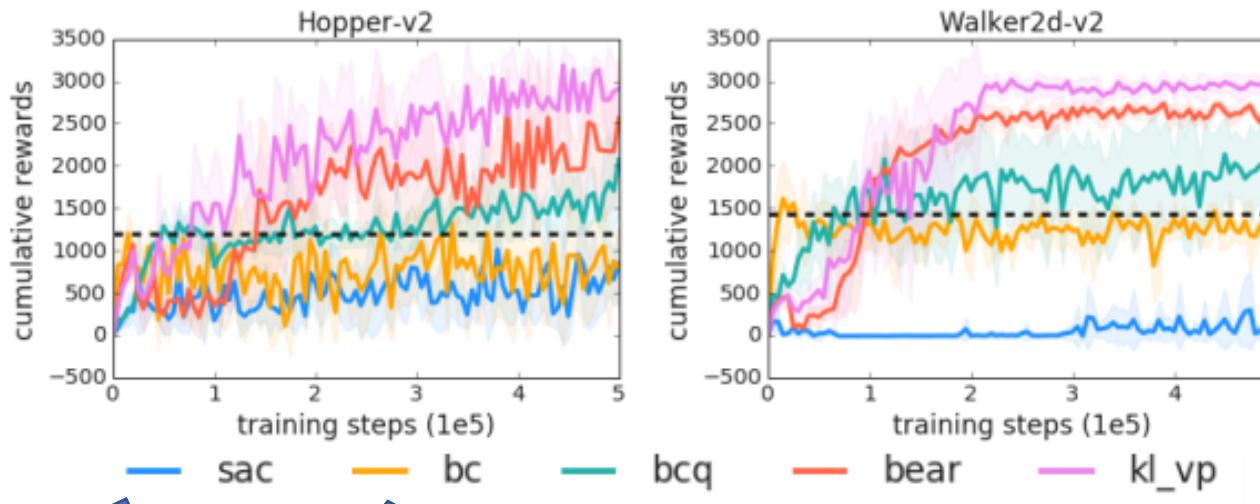
estimate $w(s, a) = \frac{d^{\pi_{\theta}}(s, a)}{d^{\pi_{\beta}}(s, a)}$ (implicitly use this for RL with a modified objective)

$$d^{\pi_{\beta}}(s', a')w(s', a') = (1-\gamma)p_0(s')\pi_{\theta}(a'|s') + \gamma \sum_{s,a} \pi_{\theta}(a'|s')p(s'|s, a)d^{\pi_{\beta}}(s, a)w(s, a)$$

As discussed before...

How do policy constraint methods perform?

Dataset collected from a mixture of random and “mediocre” policies



Naïve off-policy RL

Behavioral cloning

Policy constraint methods (BEAR: support, BCQ, BRAC: distribution)

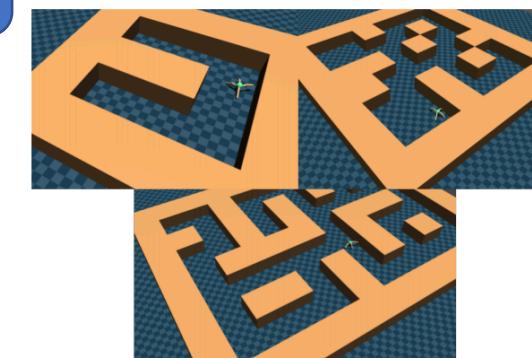
Task Name	BC	SAC	BEAR	BRAC-p
antmaze-umaze	65.0	0.0	73.0	50.0
antmaze-umaze-diverse	55.0	0.0	61.0	40.0
antmaze-medium-play	0.0	0.0	0.0	0.0
antmaze-medium-diverse	0.0	0.0	8.0	0.0
antmaze-large-play	0.0	0.0	0.0	0.0
antmaze-large-diverse	0.0	0.0	0.0	0.0

partially trained RL agent

Better than BC! (shows offline RL can do better than supervised learning)

Different choices of D matter, though not totally as expected, so other factors at play

But do not solve all tasks, especially harder ones, so work to do!



notice that kl_vp (i.e. BRAC) is a distribution matching algo and performs better than BEAR. Even though in theory, support constraint methods should be better, in practice that's not necessarily the case

Shortcomings of Policy Constraint Methods

Need estimation of the behavior policy

$$\pi_\theta := \arg \max_{\pi_\theta} \mathbb{E}_{s \sim \mathcal{D}, a \sim \pi_\theta(a|s)} [Q(s, a)] \text{ s.t. } D(\pi_\theta(a|s), \pi_\beta(a|s)) \leq \varepsilon$$

we don't know what
this is, we have to
estimate it from the
dataset

estimated from data

If the behavior policy is wrongly estimated (e.g, when we fit unimodal policy to multimodal data),
policy constraint methods can fail dramatically (e.g., AntMaze results from before)

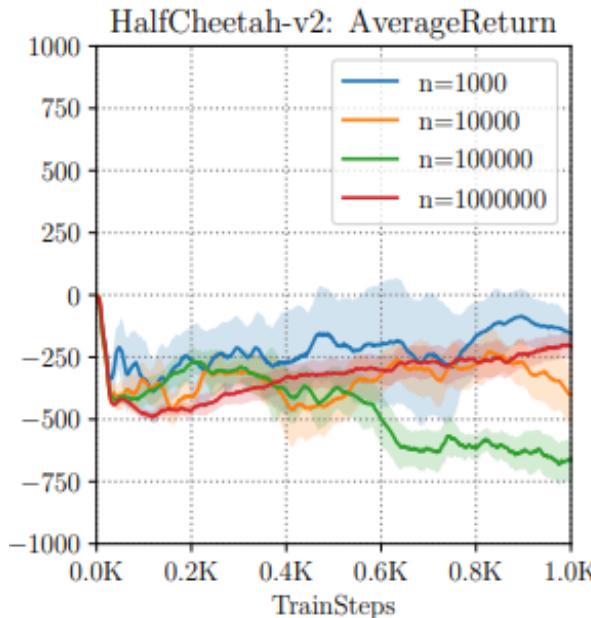
Often tend to be too conservative

If we know that a certain state has all actions with 0 reward, we do not care about
constraining the policy there, since we will not be worse...

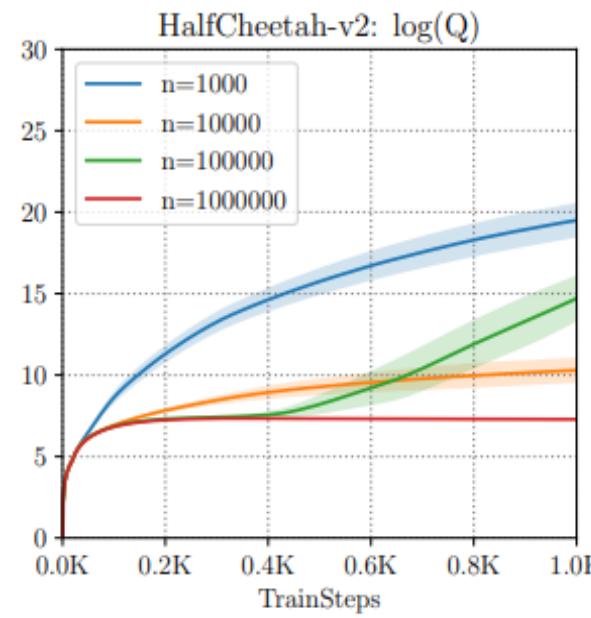
Can we do better?

Let's revisit the motivating example....

policy constraint methods seek to constrain the actions taken by the policy π_{θ} such that they are similar to the actions you'd take using π_{β} . By doing this, they indirectly tackle the overestimation issue.



how well it does



how well it *thinks*
it does (Q-values)

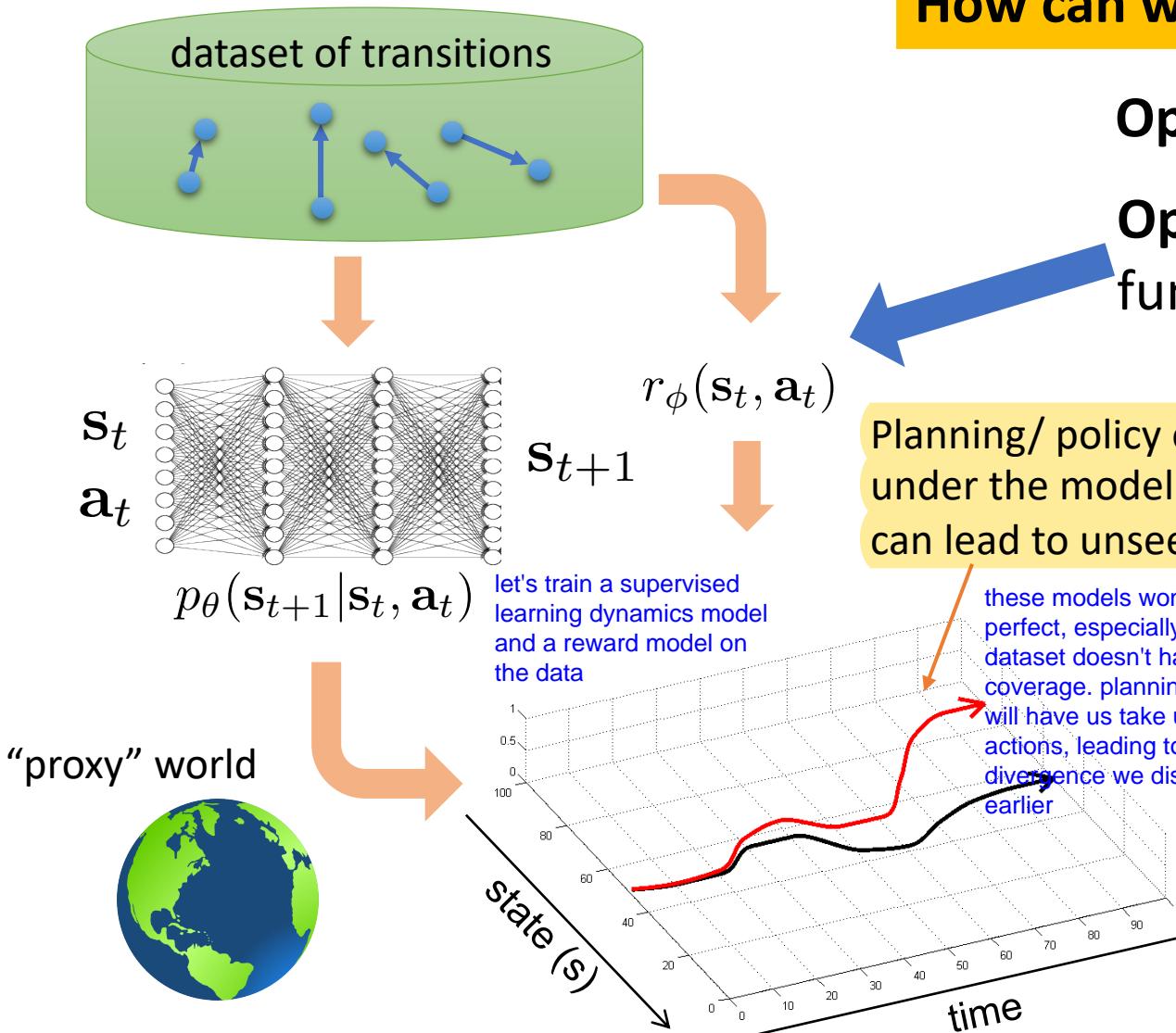
Can we directly tackle the overestimation issue?

"Conservative" Methods

Yes! Model-based and model-free methods

instead of
constraining the
action distributions,
which is not the
most optimal way to
solve this

Conservative Model-Based Offline RL



How can we avoid unseen outcomes?

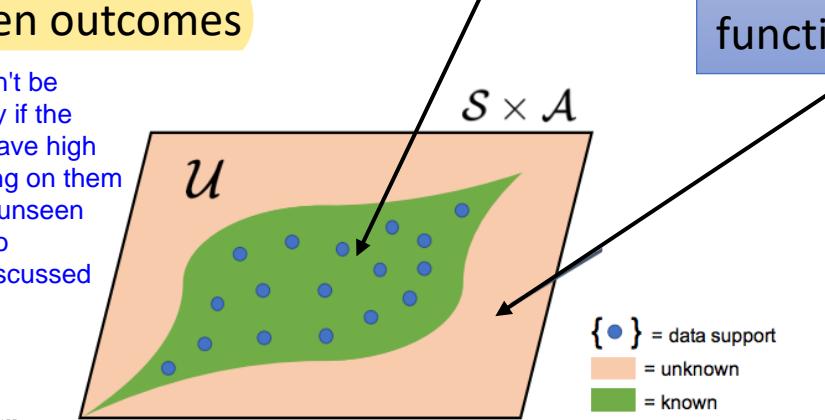
Option 1: Policy constraint (as before) as we discussed, this tends to be overly conservative

Option 2: Make the learned reward function conservative

essentially, make the reward function take into account OOD or unseen outcomes. Then run typical model-based RL using this

Keep unaltered reward function

Make reward function conservative



{●} = data support
○ = unknown
■ = known

Conservative Reward Functions in Offline RL

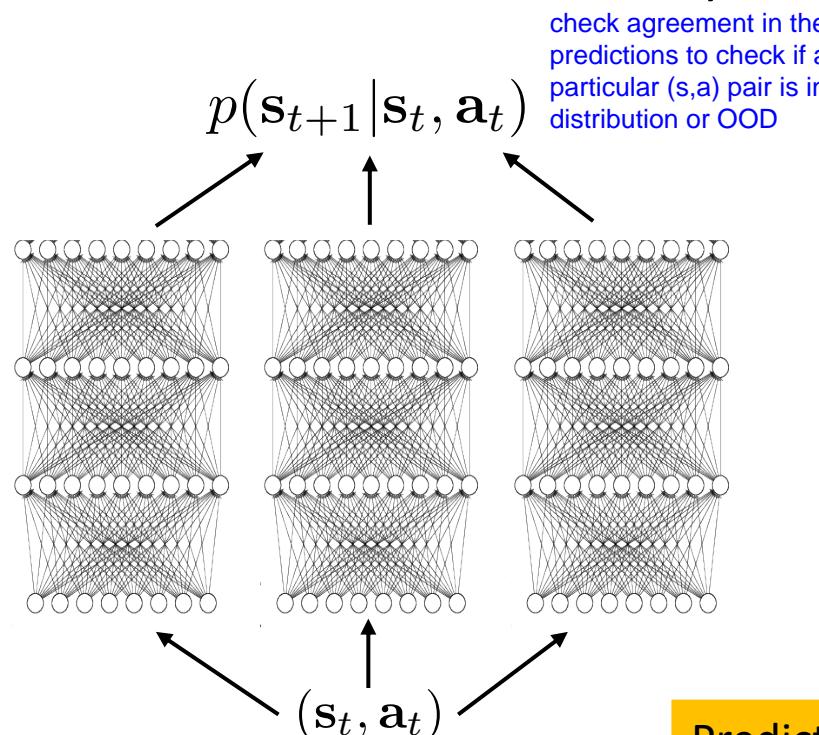
How do we obtain conservative reward functions?

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

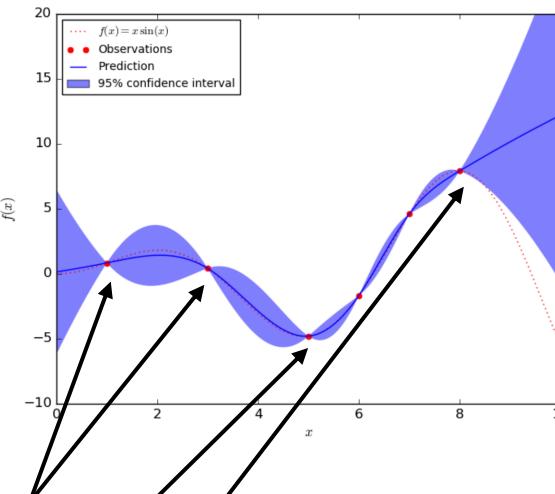
Some measure of
“uncertainty”

$u(s, a)$ provides
some indication of
whether (s, a) is
seen or unseen/
OOD

How can we estimate uncertainty in model-based offline RL?



Train an ensemble of dynamics models
and check agreement in their predictions



Predictions are certain
at observed points

Different models will make mistakes
differently on **unseen** points

the variance in predictions on unseen data points will be high

Model-Based Offline RL Algorithms

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

Option 1: Disagreement

$$dis_{i,j}(\mathbf{s}_t, \mathbf{a}_t) = \mathbb{E}_{\mathbf{s}_{t+1}^i \sim p_i(\mathbf{s}_{t+1} | \mathbf{s}_t, \mathbf{a}_t)} \mathbb{E}_{\mathbf{s}_{t+1}^j \sim p_j(\mathbf{s}_{t+1} | \mathbf{s}_t, \mathbf{a}_t)} [||s_{t+1}^i - s_{t+1}^j||]$$

if the disagreement is greater than some threshold,
subtract the maximum possible reward

some L1/L2
distance

$$u(\mathbf{s}, \mathbf{a}) = r_{\max} \quad \text{if} \quad \max_{i,j} dis_{i,j}(\mathbf{s}, \mathbf{a}) > \text{threshold}$$

MOReL (Kidambi et al. 2020)

Option 2: Prediction uncertainty

$$p_i(\mathbf{s}' | \mathbf{s}, \mathbf{a}) = \mathcal{N}(\mu_i(\mathbf{s}, \mathbf{a}), \Sigma_i(\mathbf{s}, \mathbf{a}))$$

introduce uncertainty in the
models by making
predictions based on a
gaussian distribution. Make
 $u(\mathbf{s}, \mathbf{a})$ that max frobenius
norm of the covariance
matrix

$$u(\mathbf{s}, \mathbf{a}) = \max_j \|\Sigma_j(\mathbf{s}, \mathbf{a})\|_F$$

MOPO (Yu et al. 2020)

Generate model data, run
model-free off-policy RL
algorithm (**MBPO, Dyna**)

Run trajectory optimization or
plan under the model (**Planning**)

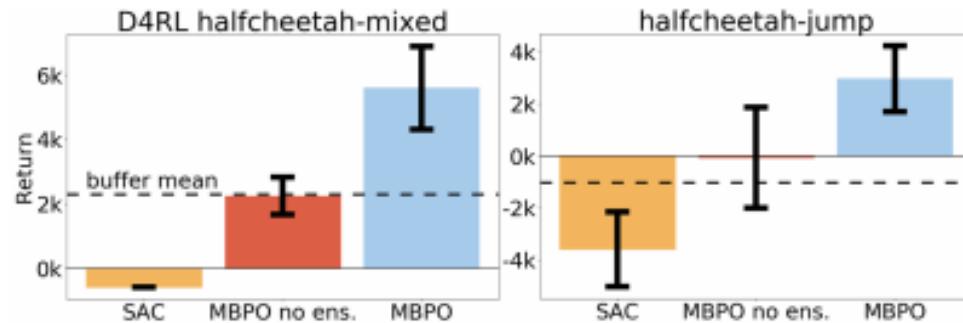
Both methods do not alter
the dynamics at all....

Should we? Or not?

now, after altering the reward
function by subtracting $u(\mathbf{s}, \mathbf{a})$, we
can either use the models in a
Dyna-style algo, or we can run a
straight planning algo. Both
options have been explored in
literature and they have trade offs.

We obtained the reward function, how do we
solve the RL problem now?

How does Model-Based Offline RL perform?



Naïve model-based methods can help with sampling error challenges if models can be learned accurately

Conservatism **does** help with model-based RL methods on “narrow” datasets

Dataset type	Environment	BC	MOPO (ours)	MBPO	SAC	BEAR	BRAC-v
random	halfcheetah	2.1	31.9 ± 2.8	30.7 ± 3.9	30.5	25.5	28.1
random	hopper	1.6	13.3 ± 1.6	4.5 ± 6.0	11.3	9.5	12.0
random	walker2d	9.8	13.0 ± 2.6	8.6 ± 8.1	4.1	6.7	0.5
medium	halfcheetah	36.1	40.2 ± 2.7	28.3 ± 22.7	-4.3	38.6	45.5
medium	hopper	29.0	26.5 ± 3.7	4.9 ± 3.3	0.8	47.6	32.3
medium	walker2d	6.6	14.0 ± 10.1	12.7 ± 7.6	0.9	33.2	81.3
mixed	halfcheetah	38.4	54.0 ± 2.6	47.3 ± 12.6	-2.4	36.2	45.9
	hopper	11.8	92.5 ± 6.3	49.8 ± 30.4	1.9	10.8	0.9
	walker2d	11.3	42.7 ± 8.3	22.2 ± 12.7	3.5	25.3	0.8
med-expert	halfcheetah	35.8	57.9 ± 24.8	9.7 ± 9.5	1.8	51.7	45.3
med-expert	hopper	111.9	51.7 ± 42.9	56.0 ± 34.5	1.6	4.0	0.8
med-expert	walker2d	6.4	55.0 ± 19.1	7.6 ± 3.7	-0.1	26.0	66.6

Generally **better** than policy constraint methods: less conservative in many cases!

Value Function Regularization for Offline RL

So far, we have looked at handling overestimation via penalizing the reward function

Can we penalize value functions or Q-functions directly?

Way 1: Add penalty based on policy constraint to the value function

$$Q(s, a) \leftarrow r(s, a) + \gamma \mathbb{E}_{a' \sim \pi_\theta(a'|s')} [Q(s', a')] - \alpha D(\pi_\theta, \pi_\beta)$$

But this still uses a policy constraint.....

KL-control (Kakade '02, Fox et al. '16, Haarnoja et al. '17)

BRAC-v (Wu et al. 2019), KL-control (Jaques et al. 2019)

Why can this work?

Way 2: Change the objective to make this behavior automatically kick in!

TD error will push values on **seen** state-action pairs up, while we minimize others

Minimize the “big” Q-values

Standard TD error objective

$$\min_Q \max_\mu \mathbb{E}_{s \sim \mathcal{D}} \mathbb{E}_{a \sim \mu(a|s)} [Q(s, a)] + \frac{1}{2\alpha} \mathbb{E}_{s, a, s' \sim \mathcal{D}} [(Q(s, a) - y(s, a))^2]$$

additional regularizer that serves to minimize the "big" Q-values. So it identifies actions with high Q-values, and minimizes and pulls those Q-values down

$$y(s, a) = r(s, a) + \gamma \mathbb{E}_{a' \sim \pi_\theta(a'|s')} [Q(s', a')]$$

Target values

essentially, the TD error pushes up Q-values on seen state-action pairs. This happens through the bellman backup. We then pull the Q-values down for unseen actions where the Q-values are high. So you make sure the Q-values on actions we haven't seen are not overestimated.

Conservative Q-Learning

CQL-v1

$$Q_{\text{CQL}}(\mathbf{s}, \mathbf{a}) := \arg \min_Q \max_\mu \mathbb{E}_{\mathbf{s} \sim \mathcal{D}} \mathbb{E}_{\mathbf{a} \sim \mu(\mathbf{a}|\mathbf{s})} [Q(\mathbf{s}, \mathbf{a})] + \frac{1}{2\alpha} \mathbb{E}_{\mathbf{s}, \mathbf{a}, \mathbf{s}' \sim \mathcal{D}} [(Q(\mathbf{s}, \mathbf{a}) - y(\mathbf{s}, \mathbf{a}))^2]$$

$$\forall \mathbf{s} \in \mathcal{D}, \mathbf{a}, \quad Q_{\text{CQL}}(\mathbf{s}, \mathbf{a}) \leq Q(\mathbf{s}, \mathbf{a})$$

- + Conservative (so no overestimation)
- ✗ Can be too conservative!

CQL-v2

Minimize the big Q-values

$$Q_{\text{CQL}}(\mathbf{s}, \mathbf{a}) := \arg \min_Q \max_\mu \left(\underbrace{\mathbb{E}_{\mathbf{s} \sim \mathcal{D}} \mathbb{E}_{\mathbf{a} \sim \mu(\mathbf{a}|\mathbf{s})} [Q(\mathbf{s}, \mathbf{a})]}_{\text{Minimize the big Q-values}} - \underbrace{\mathbb{E}_{\mathbf{s}, \mathbf{a} \sim \mathcal{D}} [Q(\mathbf{s}, \mathbf{a})]}_{\text{Maximize the data Q-values}} \right) + \frac{1}{2\alpha} \mathbb{E}_{\mathbf{s}, \mathbf{a}, \mathbf{s}' \sim \mathcal{D}} [(Q(\mathbf{s}, \mathbf{a}) - y(\mathbf{s}, \mathbf{a}))^2]$$

$$\forall \mathbf{s} \in \mathcal{D}, \mathbf{a}, \quad Q_{\text{CQL}}(\mathbf{s}, \mathbf{a}) \leq Q(\mathbf{s}, \mathbf{a})$$



Maximize the data Q-values

Value function is still a lower-bound,
i.e. policy value is lower-bounded,

$$\forall \mathbf{s} \in \mathcal{D}, \quad V_{\text{CQL}}(\mathbf{s}) \leq V(\mathbf{s})$$



Much “tighter”

CQL Algorithm:

1. Learn \hat{Q}_{CQL}^π using offline data \mathcal{D} .
2. Optimize policy w.r.t. \hat{Q}_{CQL}^π : $\pi \leftarrow \arg \max_\pi \mathbb{E}_\pi [\hat{Q}_{\text{CQL}}^\pi]$.

Related ideas: Piot et al. ‘14, Hester et al. ‘17: use “margin loss” instead of Q-function difference

Value Regularization in Practice

Minimize "soft"-maximum Q-values and maximize data Q-values

$$\min_Q \alpha \mathbb{E}_{\mathbf{s} \sim \mathcal{D}} \left[\log \sum_{\mathbf{a}} \exp(Q(\mathbf{s}, \mathbf{a})) - \mathbb{E}_{\mathbf{a} \sim \hat{\pi}_\beta(\mathbf{a}|\mathbf{s})} [Q(\mathbf{s}, \mathbf{a})] \right] + \frac{1}{2} \mathbb{E}_{\mathbf{s}, \mathbf{a}, \mathbf{s}' \sim \mathcal{D}} \left[(Q - \hat{B}^{\pi_k} \hat{Q}^k)^2 \right].$$

Policy constraint penalty vs. regularized objectives

$$Q(\mathbf{s}, \mathbf{a}) \leftarrow r(\mathbf{s}, \mathbf{a}) + \gamma \mathbb{E}_{\mathbf{a}' \sim \pi_\theta} [Q(\mathbf{s}', \mathbf{a}')] - \alpha D(\pi_\theta, \pi_\beta)$$

- 
1. update Q_ϕ to decrease ◇
 2. update π_θ to increase $E_{\mathbf{s} \sim d^{\pi_\theta}(\mathbf{s}), \mathbf{a} \sim \pi_\theta(\mathbf{a}|\mathbf{s})} [Q_\phi(\mathbf{s}, \mathbf{a})]$

A different penalty, but very similar

With tabular Q-functions

$$Q(\mathbf{s}, \mathbf{a}) \leftarrow r(\mathbf{s}, \mathbf{a}) + \gamma \mathbb{E}_{\mathbf{a}' \sim \pi_\theta} [Q(\mathbf{s}', \mathbf{a}')] - \alpha D_{CQL}(\pi_\theta, \pi_\beta)$$

With Q-function approximators (in practice)

Value Regularization = "Q-function aware" penalty

"projection" operator, a concise notation to represent function approximators

$$Q(\mathbf{s}, \mathbf{a}) \leftarrow \prod_\phi (r(\mathbf{s}, \mathbf{a}) + \gamma \mathbb{E}_{\mathbf{a}' \sim \pi_\theta} [Q(\mathbf{s}', \mathbf{a}')] - \alpha D_{CQL}(\pi_\theta, \pi_\beta)) \times \cancel{\text{X}}$$

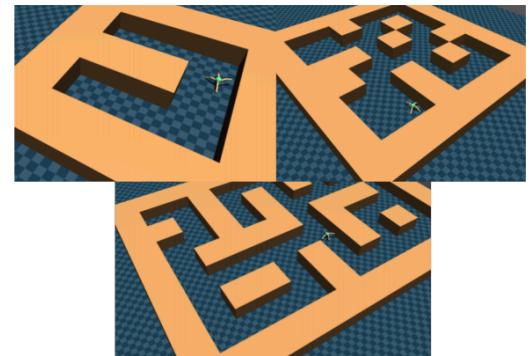
$$Q(\mathbf{s}, \mathbf{a}) \leftarrow \prod_\phi (r(\mathbf{s}, \mathbf{a}) + \gamma \mathbb{E}_{\mathbf{a}' \sim \pi_\theta} [Q(\mathbf{s}', \mathbf{a}')] - \alpha f(\pi_\theta, \pi_\beta, Q_\phi))$$

How does CQL perform?

Learned policy value - Actual policy value

CQL learns conservative Q-values

Task Name	CQL(\mathcal{H})	CQL (Eqn. 1)	Ensemble(2)	Ens.(4)	Ens.(10)	Ens.(20)	BEAR
hopper-medium-expert	-43.20	-151.36	3.71e6	2.93e6	0.32e6	24.05e3	65.93
hopper-mixed	-10.93	-22.87	15.00e6	59.93e3	8.92e3	2.47e3	1399.46
hopper-medium	-7.48	-156.70	26.03e12	437.57e6	1.12e12	885e3	4.32



Imitation

Naïve off-policy RL

Policy constraints

“Way 1”

CQL variants

Domain	Task Name	BC	SAC	BEAR	BRAC-p	BRAC-v	CQL(\mathcal{H})	CQL(ρ)
AntMaze	antmaze-umaze	65.0	0.0	73.0	50.0	70.0	74.0	73.5
	antmaze-umaze-diverse	55.0	0.0	61.0	40.0	70.0	84.0	61.0
	antmaze-medium-play	0.0	0.0	0.0	0.0	0.0	61.2	4.6
	antmaze-medium-diverse	0.0	0.0	8.0	0.0	0.0	53.7	5.1
	antmaze-large-play	0.0	0.0	0.0	0.0	0.0	15.8	3.2
	antmaze-large-diverse	0.0	0.0	0.0	0.0	0.0	14.9	2.3

CQL performs better due to stitching!

so far we've seen three approaches:
 1) policy constraint methods, 2)
 Model-based methods, 3) Model-free
 methods.

Uncertainty-Based Offline RL Methods

But what if our model or Q-function can generalize on unseen (state, action) pairs? In that case, we wouldn't necessarily want to penalize them and be conservative! But figuring out whether or not we can generalize is tricky.

Promise: Can do better by taking into account **when** our Q-function/model can generalize!

two-sides of the spectrum: A) traditional off-policy, B) add conservatism. By understanding how well we can generalize, we can modulate where we are between the two extremes.

Model-based methods

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

We already do it! – ensembles of models, disagreement, etc

model-based methods fall into the uncertainty-based category.

Model-free methods

Computing an uncertainty function and running explicit RL on it tends to be not be much different than policy constraint methods



How can we model uncertainty sets over Q-functions?

A different approach:
 directly model Q-function uncertainty

$\mathcal{P}_{\mathcal{D}}(Q^{\pi})$ = uncertainty set of Q-functions given the dataset

$$\pi_{k+1} \leftarrow \arg \max_{\pi} \mathbb{E}_{s \sim \mathcal{D}} \left[\mathbb{E}_{a \sim \pi(a|s)} \left[\mathbb{E}_{Q_{k+1}^{\pi} \sim \mathcal{P}_{\mathcal{D}}(Q^{\pi})} [Q_{k+1}^{\pi}(s, a)] - \alpha \text{Unc}(\mathcal{P}_{\mathcal{D}}(Q^{\pi})) \right] \right]$$

conservative estimate

expected Q-value (standard)

conservative when uncertain

Uncertainty-Based Offline RL Methods

Bootstrap ensembles of Q-functions

average dirac delta distribution over all possible q-functions

$$\mathcal{P}_{\mathcal{D}}(Q^{\pi}) \approx \frac{1}{N} \sum_{j=1}^N \delta[Q^{\pi} = Q_{\phi_j}]$$

$$\pi_{k+1} \leftarrow \arg \max_{\pi} \mathbb{E}_{\mathbf{s} \sim \mathcal{D}} \left[\underbrace{\mathbb{E}_{\mathbf{a} \sim \pi(\mathbf{a}|\mathbf{s})} \left[\mathbb{E}_{Q_{k+1}^{\pi} \sim \mathcal{P}_{\mathcal{D}}(Q^{\pi})} [Q_{k+1}^{\pi}(\mathbf{s}, \mathbf{a})] - \alpha \text{Unc}(\mathcal{P}_{\mathcal{D}}(Q^{\pi})) \right]}_{\text{conservative estimate}} \right]$$

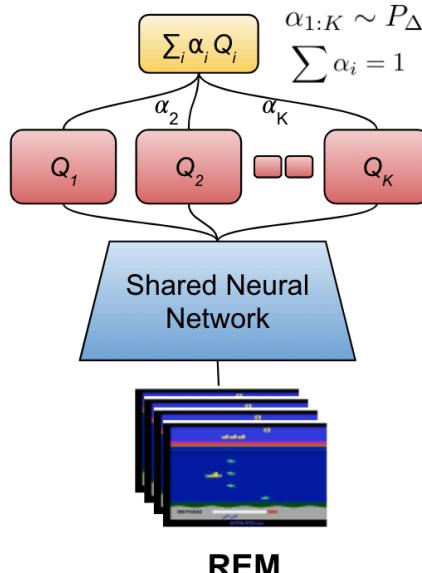
N Q-functions $Q_{\phi_1}, Q_{\phi_2}, \dots, Q_{\phi_N}$

Train each Q-function independently

$$\min_{\phi_i} \mathbb{E}_{\mathbf{s}, \mathbf{a} \sim \mathcal{D}} [(Q_{\phi_i}(\mathbf{s}, \mathbf{a}) - y_i(\mathbf{s}, \mathbf{a}))^2]$$

so we underestimate uncertainty and has shown not to work well in practice so far

Very little diversity in an ensemble and thus incorrect uncertainty (Fujimoto et al. 2019, Kumar et al. 2019)



Random Ensemble Mixture (REM)

$$Q(\mathbf{s}, \mathbf{a}) = \sum_{j=1}^N \alpha_j Q_{\phi_j}(\mathbf{s}, \mathbf{a})$$

$$\text{Unc}(Q_{\phi_1}, \dots, Q_{\phi_N}) = 0$$

...but weights randomly chosen!

Overall, a general solution that works well in each case has been hard to obtain

this approach has been shown to work. they didn't explain this well

Works well with high-coverage datasets

we partition our dataset into independent partitions (i.e. bootstrapping). We make sure that each Q-function has their own, independent target function too.

Summary of Offline RL Methods

- **Policy constraint** often work well when the behavior policy distribution is “easy” to model, but tend to be more conservative
- **Model-based offline RL** often work quite well when the data distribution is relatively high-coverage, and it is easy to learn a model you can approximate the dynamics and reward function with high accuracy
- **Value-regularized offline RL** methods tend to be less conservative and can be augmented with either models or policy constraints, if needed

Method Type	Modifies	Modification Type	Extra requirements
Policy constraints	Learned policy	Probabilistic constraint, support constraint, state-marginal constraints	Behavior policy estimation (except implicit constraints)
Model-based	Reward function	Subtracts an uncertainty function	Defining an uncertainty/conservatism metric
Value regularization	Q-function	Modified Q-function training objective	Defining a regularizer
Uncertainty-based	Q-function or model	Subtract an uncertainty (one approach)	Uncertainty estimation

Additional readings: Modern Offline RL Methods

- **Summary/Survey:** Levine, Kumar, Tucker, Fu (2020). Offline Reinforcement Learning: Tutorial, Survey and Perspectives on Open Problems. <https://arxiv.org/abs/2005.01643>
 - **Policy Constraint Algorithms**
 - Fujimoto et al. (2019) Off-Policy Reinforcement Learning without Exploration. ICML 2019
 - Laroche et al. (2019) Safe Policy Improvement via Baseline Bootstrapping. ICML 2019.
 - Nadjahi et al. (2019). Safe Policy Improvement with Soft Baseline Bootstrapping. ECML-PKDD 2019.
 - - Simao et al. (2019). Safe Policy Improvement with an Estimated Behavior Policy.
 - Kumar et al. (2019) Stabilizing Off-Policy Reinforcement Learning via Bootstrapping Error Reduction. NeurIPS 2019
 - Wu et al. (2019). Behavior Regularized Offline Reinforcement Learning.
 - Jaques et al. (2019). Way Off-Policy Batch Deep Reinforcement Learning of Implicit Human Preferences in Dialog.
 - Peng et al. (2019). Advantage-Weighted Regression: Simple and Scalable Off-Policy Reinforcement Learning.
 - Siegel et al. (2019). Keep Doing What Worked: Behavior Modelling Priors for Offline RL. ICLR 2020.
 - Ghasemipour et al. (2020). EMaQ: Expected-Max Q-Learning Operator for Simple Yet Effective Offline and Online RL.
 - Liu et al. (2020). Provably Good Batch RL without Great Exploration. NeurIPS 2020.
 - Wang et al. (2020). Critic-regularized Regression. NeurIPS 2020.
 - Chen et al. (2020). BAIL: Best-Action Imitation Learning for Batch Deep Reinforcement Learning. NeurIPS 2020.
 - **Model-Based [MB]/Value Regularization [V]/Uncertainty-aware [U] algorithms**
 - [MB] Yu et al. (2020). MOPO: Model-Based Offline Policy Optimization. NeurIPS 2020.
 - [MB] Kidambi et al. (2020). Model-Based Offline Reinforcement Learning. NeurIPS 2020.
 - [V] Kumar et al. (2020). Conservative Q-Learning for Offline Reinforcement Learning. NeurIPS 2020.
 - [U] Agarwal et al. (2020). An Optimistic Perspective on Offline Reinforcement Learning. ICML 2020.
 - [U] Sonabend-W et al. (2020). Expert-Supervised RL for Offline Policy Learning and Evaluation. NeurIPS 2020.
- ...and many more

1. Reinforcement learning primer
2. The offline reinforcement learning problem
3. Classic batch reinforcement learning algorithms
 - i. Importance-sampled policy gradients
 - ii. Least-squares temporal difference (LSTD)
4. Modern offline reinforcement learning algorithms
 - i. Policy constraint methods
 - ii. Model-based offline RL
 - iii. Value function regularization
 - iv. Uncertainty-based methods
 - v. Comparisons and discussion
5. **Off-policy evaluation and model selection**
 - i. Problem statement
 - ii. Direct methods and importance sampling methods
 - iii. Doubly robust methods
6. How should we evaluate offline RL methods?
 - i. When and why does offline RL work?
 - ii. How do we benchmark offline RL methods?
7. Conclusions, discussion, and open problems

Off-Policy Evaluation and Model Selection

Problem Statement: Rather than returning a good policy, find me the value of a **given policy**, **without running this policy in the environment**

$$\pi \xrightarrow{\mathcal{D}} V^\pi(s)$$

$$V^{\pi_1}(s) > V^{\pi_2}(s)?$$

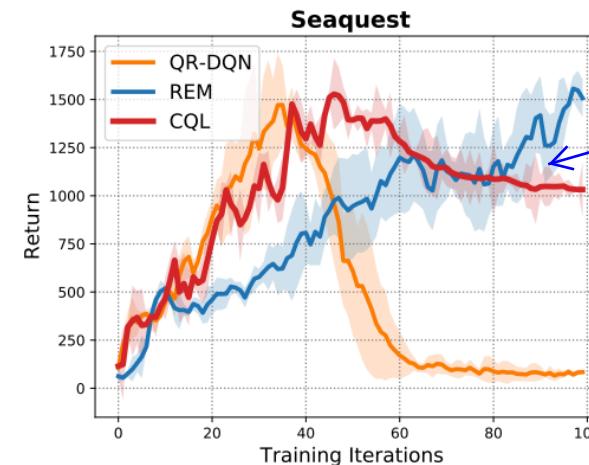
What is the use of OPE in offline RL?

Model selection: selecting hyperparameters/policies

Why do we need model selection in offline RL?

Similar to supervised learning methods, excessive training on the same offline dataset can produce poor solutions. Moreover, it is unclear how to select hyperparameters.

offline algorithms can be extremely sensitive to hyperparameters, and can be the difference between a method working and not working!

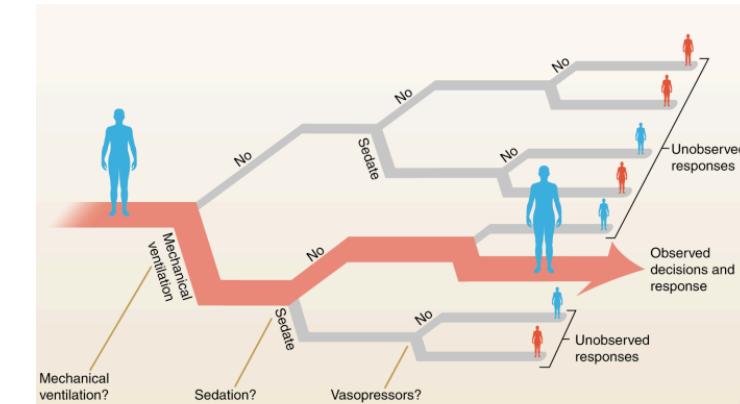


NeurIPS 2020

Learning. '2019

Wu, Tucker, Nachum. Behavior Regul

Agarwal, Schuurmans, Norouzi. An Optimalistic Perspective on Offline Reinforcement Learning. ICML 2020



(Figure from Gottesman et al. Nature Medicine, 2019)

notice how performance drops after a certain number of training iterations

Direct Methods for Off-Policy Evaluation

- **Model-based approach:** Fit a transition and reward function model. Compute return of the policy under the model.

$$\hat{p}(\mathbf{s}_{t+1}|\mathbf{s}_t, \mathbf{a}_t) \quad \hat{r}(\mathbf{s}_t, \mathbf{a}_t) \quad \hat{J}(\pi) = \mathbb{E}_{\mathbf{s}_0, \mathbf{s}_{t+1} \sim \hat{p}, \mathbf{a}_t \sim \pi} \left[\sum_t \gamma^t \hat{r}(\mathbf{s}_t, \mathbf{a}_t) \right]$$

policy to evaluate

- **Fitted-Q Evaluation (FQE):** Evaluate a policy by training a Q-function using Bellman consistency for the chosen policy

$$Q^\pi(\mathbf{s}, \mathbf{a}) = r(\mathbf{s}, \mathbf{a}) + \gamma \mathbb{E}_{\mathbf{s}' \sim p(\mathbf{s}'|\mathbf{s}, \mathbf{a}), \mathbf{a}' \sim \pi(\mathbf{a}'|\mathbf{s}')} [Q^\pi(\mathbf{s}', \mathbf{a}')]$$

$$\hat{Q}_\phi := \arg \min_{\phi} \mathbb{E}_{\mathbf{s}, \mathbf{a}, \mathbf{s}' \sim \mathcal{D}} [(Q_\phi(\mathbf{s}, \mathbf{a}) - (r(\mathbf{s}, \mathbf{a}) + \gamma \mathbb{E}_{\mathbf{a}' \sim \pi} [Q_\phi(\mathbf{s}', \mathbf{a}')])))^2]$$

Works quite well in practice

$$\hat{J}(\pi) = \mathbb{E}_{\mathbf{s}_0, \mathbf{a}_0 \sim \pi(\mathbf{a}_0|\mathbf{s}_0)} [\hat{Q}_\phi(\mathbf{s}_0, \mathbf{a}_0)]$$

Importance Sampling OPE Methods

- Per-decision importance sampling/ Inverse propensity scoring (IPS): Compute importance ratios and re-weight trajectories

$$\begin{aligned}
 J(\pi_\theta) &= \mathbb{E}_{\tau \sim \pi_\beta(\tau)} \left[\frac{\pi_\theta(\tau)}{\pi_\beta(\tau)} \sum_{t=0}^H \gamma^t r(\mathbf{s}, \mathbf{a}) \right] \\
 &= \mathbb{E}_{\tau \sim \pi_\beta(\tau)} \left[\left(\prod_{t=0}^H \frac{\pi_\theta(\mathbf{a}_t | \mathbf{s}_t)}{\pi_\beta(\mathbf{a}_t | \mathbf{s}_t)} \right) \sum_{t=0}^H \gamma^t r(\mathbf{s}, \mathbf{a}) \right] \approx \sum_{i=1}^n w_H^i \sum_{t=0}^H \gamma^t r_t^i,
 \end{aligned}$$

High variance

Sum over the dataset

Markovian structure?

- Step-wise importance sampling estimator: Can replace with step-wise weights, rather than trajectory-level weights

$$J(\pi_\theta) \approx \frac{1}{|\mathcal{D}|} \sum_{i=1}^n \sum_{t=0}^H \gamma^t \left(\prod_{j=0}^t \frac{\pi_\theta(\mathbf{a}_j | \mathbf{s}_j)}{\pi_\beta(\mathbf{a}_j | \mathbf{s}_j)} \right) r_t(\mathbf{s}_t, \mathbf{a}_t)$$

Improves over trajectory level weights

"surrogate" for state frequency

Doubly-Robust Importance Sampling

Let's add a learned Q-function to reduce variance, which still retaining the unbiased nature

Need to evaluate $V(\mathbf{s}) := \mathbb{E}_{\mathbf{a} \sim \pi(\mathbf{a}|\mathbf{s})}[r(\mathbf{s}, \mathbf{a})]$

$$\hat{V}(\mathbf{s}) = \mathbb{E}_{\mathbf{a} \sim \pi(\mathbf{a}|\mathbf{s})}[\hat{r}(\mathbf{s}, \mathbf{a})]$$

Reward model,
separate estimator

Importance sampling

$$\hat{V}_{\text{IS}}(\mathbf{s}) := \frac{\pi(\mathbf{a}|\mathbf{s})}{\pi_{\beta}(\mathbf{a}|\mathbf{s})} [r(\mathbf{s}, \mathbf{a})]$$

Importance weights + a different estimator

$$\hat{V}_{\text{DR}}(\mathbf{s}) = \hat{V}(\mathbf{s}) + \frac{\pi(\mathbf{a}|\mathbf{s})}{\pi_{\beta}(\mathbf{a}|\mathbf{s})} (r(\mathbf{s}, \mathbf{a}) - \hat{r}(\mathbf{s}, \mathbf{a}))$$

Unbiased in expectation

Take step-wise IS estimator

$$V_{\text{step-IS}}^{\pi}(\mathbf{s}_t) = \frac{\pi(\mathbf{a}_t|\mathbf{s}_t)}{\pi_{\beta}(\mathbf{a}_t|\mathbf{s}_t)} (r(\mathbf{s}_t, \mathbf{a}_t) + \gamma V_{\text{step-IS}}^{\pi}(\mathbf{s}_{t+1}))$$

$$V_{\text{DR}}^{\pi}(\mathbf{s}_t) = \hat{V}(\mathbf{s}_t) + \frac{\pi(\mathbf{a}_t|\mathbf{s}_t)}{\pi_{\beta}(\mathbf{a}_t|\mathbf{s}_t)} \left(r(\mathbf{s}_t, \mathbf{a}_t) + \gamma V_{\text{DR}}^{\pi}(\mathbf{s}_{t+1}) - \hat{Q}(\mathbf{s}_t, \mathbf{a}_t) \right)$$

Lower variance

Marginalized Importance Sampling

Importance sampling over a trajectory suffers from high variance. Can we do better via state distributions?

Main idea: instead of using $\prod_t \frac{\pi_\theta(\mathbf{a}_t|\mathbf{s}_t)}{\pi_\beta(\mathbf{a}_t|\mathbf{s}_t)}$, estimate $w(\mathbf{s}, \mathbf{a}) = \frac{d^{\pi_\theta}(\mathbf{s}, \mathbf{a})}{d^{\pi_\beta}(\mathbf{s}, \mathbf{a})}$

$$J(\pi_\theta) \approx \frac{1}{N} \sum_i w(\mathbf{s}_i, \mathbf{a}_i) r(\mathbf{s}_i, \mathbf{a}_i)$$

DualDICE: importance weights are solutions to:

$$w(\mathbf{s}, \mathbf{a}) = \arg \min_{x: \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}} \frac{1}{2} \underbrace{\mathbb{E}_{\mathbf{s}, \mathbf{a}, \mathbf{s}' \sim \mathcal{D}} [x(\mathbf{s}, \mathbf{a})^2] - \mathbb{E}_{\mathbf{s} \sim d^\pi(\mathbf{s}), \mathbf{a} \sim \pi(\mathbf{a}|\mathbf{s})} [x(\mathbf{s}, \mathbf{a})]}_{(Change\ of\ variables)} \underbrace{(\hat{\mathcal{B}}^\pi \nu)(\mathbf{s}, \mathbf{a})}_{x(\mathbf{s}, \mathbf{a}) = \nu(\mathbf{s}, \mathbf{a}) - \mathbb{E}_{\mathbf{s}' \sim p(\mathbf{s}'|\mathbf{s}, \mathbf{a}), \mathbf{a}' \sim \pi(\mathbf{a}'|\mathbf{s}')} [\nu(\mathbf{s}', \mathbf{a}')]}$$
$$\min_{\nu: \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}} \frac{1}{2} \mathbb{E}_{\mathbf{s}, \mathbf{a}, \mathbf{s}' \sim \mathcal{D}} \left[(\nu(\mathbf{s}, \mathbf{a}) - \tilde{\mathcal{B}}^\pi \nu(\mathbf{s}, \mathbf{a}))^2 \right] - \mathbb{E}_{\mathbf{s}_0 \sim d_0(\mathbf{s}_0), \mathbf{a} \sim \pi(\mathbf{a}|\mathbf{s}_0)} [\nu(\mathbf{s}_0, \mathbf{a})]$$

Very much like a Q-function

Summary of OPE and Model Selection

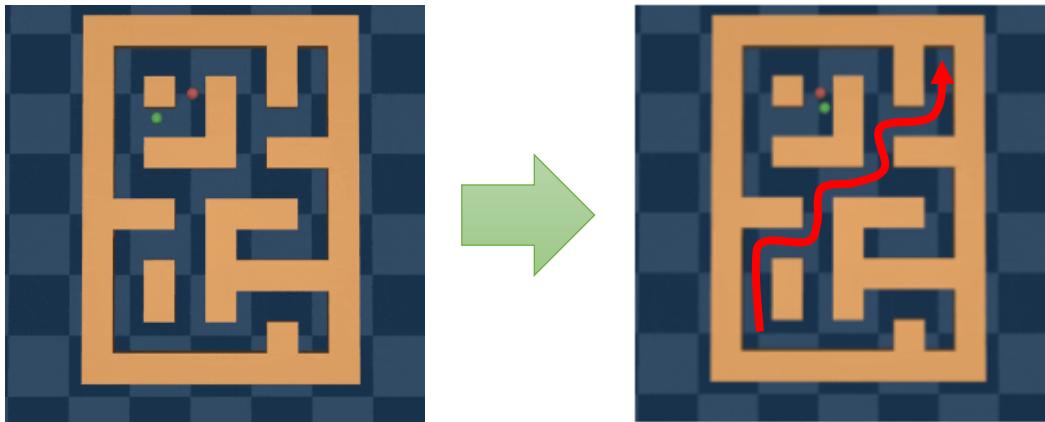
- **Direct Methods (e.g., FQE)** have been shown to work well, but potentially suffer from distributional shift similar to policy learning methods
- **Marginalized importance sampling (MIS)** methods enjoy good theoretical guarantees, however they tend to be a bit challenging in practice with function approximation (saddle point problems)
- **Importance sampling (IS)** methods suffer from high variance and produce incorrect estimates in many cases, especially when the behavior policy needs to be estimated.

(See Voloshin et al. (2019). Empirical Analysis of Off-Policy Policy Evaluation for Reinforcement Learning)

1. Reinforcement learning primer
2. The offline reinforcement learning problem
3. Classic batch reinforcement learning algorithms
 - i. Importance-sampled policy gradients
 - ii. Least-squares temporal difference (LSTD)
4. Modern offline reinforcement learning algorithms
 - i. Policy constraint methods
 - ii. Model-based offline RL
 - iii. Value function regularization
 - iv. Uncertainty-based methods
 - v. Comparisons and discussion
5. Off-policy evaluation and model selection
 - i. Problem statement
 - ii. Direct methods and importance sampling methods
 - iii. Doubly robust methods
6. **How should we evaluate offline RL methods?**
 - i. When and why does offline RL work?
 - ii. How do we benchmark offline RL methods?
7. Conclusions, discussion, and open problems

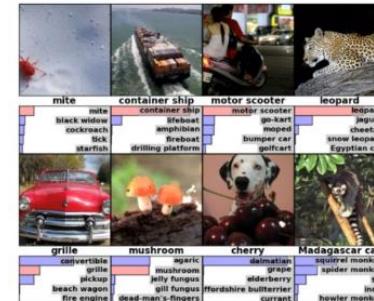
When and why does offline RL work?

1. Find the “good stuff” in a dataset full of good and bad behaviors
2. Generalization: good behavior in one place may suggest good behavior in another place
3. “Stitching”: parts of good behaviors can be recombined

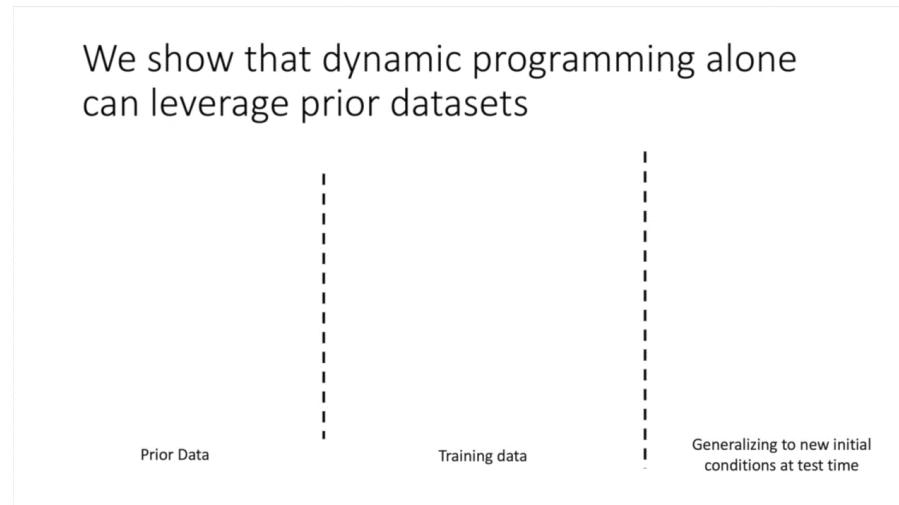


What about stochastic settings?

Just taking the best seen action is not always best!



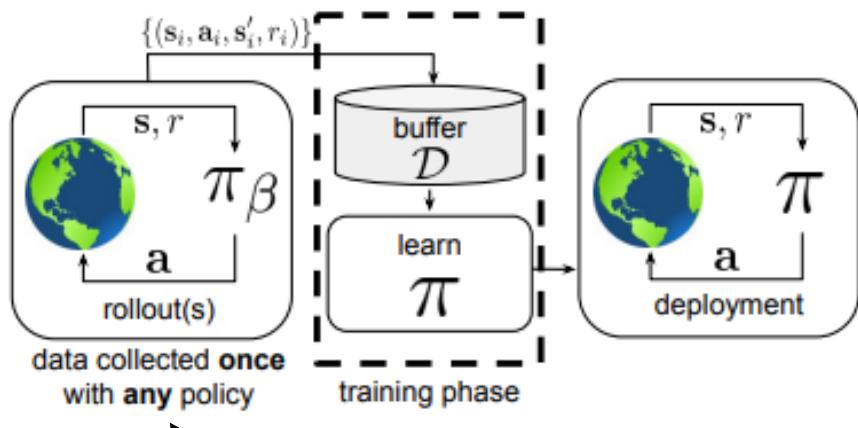
We show that dynamic programming alone can leverage prior datasets



How do we know if it's going to work?

- **Major challenge:** if we can't roll out the policy, how do we know if it will work?
- This is a huge barrier to practical use of offline RL methods
- There is no single great solution, though off-policy evaluation (OPE) methods can help
- Usually this means that when we evaluate offline RL methods, we would still use a simulator/rollouts to see if our method works
- No great answer for hyperparameter tuning just yet, but active area of research (see, e.g., the OPE section before)

How do we evaluate offline RL methods?



typical protocol in prior work:

1. train π_β with *online* RL
2. either collect data throughout training
OR
2. collect data from final policy π_β

this is not good

- In the real world, data might come from non-Markovian “policies”
 - Human users
 - Hand-engineered policies
- Must use data that is **representative of real-world settings** and **leaves lots of room for improvement**
data must be suboptimal, or else the best you can do is behavioral cloning
- Offline RL **must learn policies that are much better than the behavior policy!**
otherwise there's no point in doing offline RL

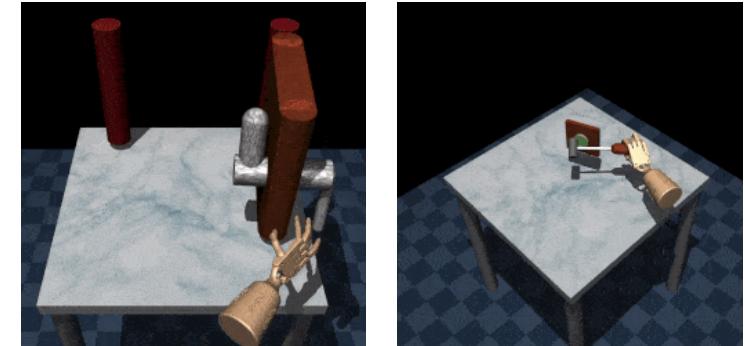
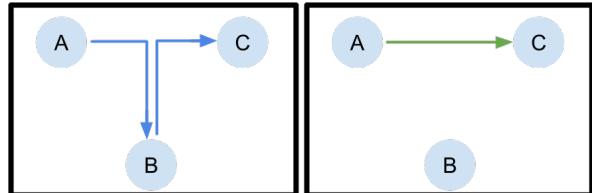
without testing these properties, we **cannot** trust that our algorithms are good!

D4RL: Datasets for Data-Driven Deep RL

What are some important principles to keep in mind?

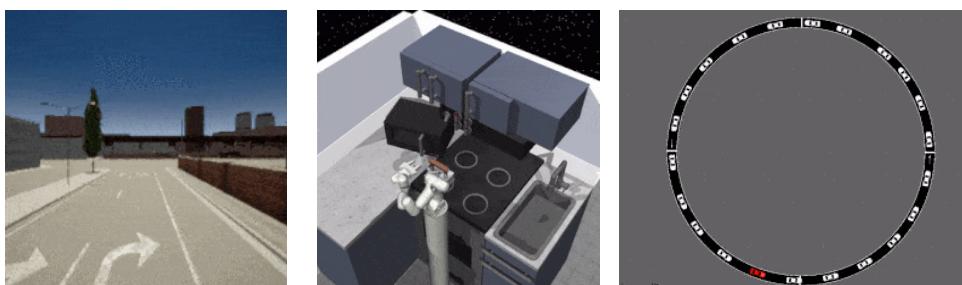
Data from non-RL policies, including data from humans

Stitching: data where dynamic programming can find much better solutions



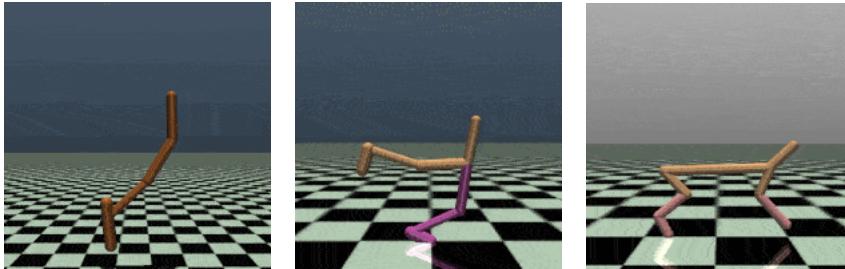
simulation & human data from Rajeswaran et al.

Realistic tasks



The D4RL Tasks

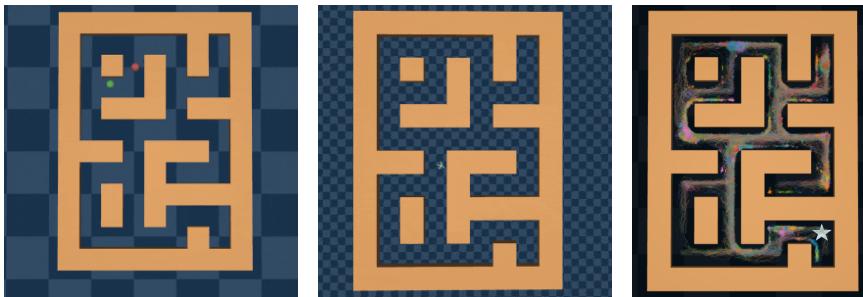
standard MuJoCo Gym-style benchmarks



standardized “expert,” “medium,” and “random” datasets introduced by Kumar et al.

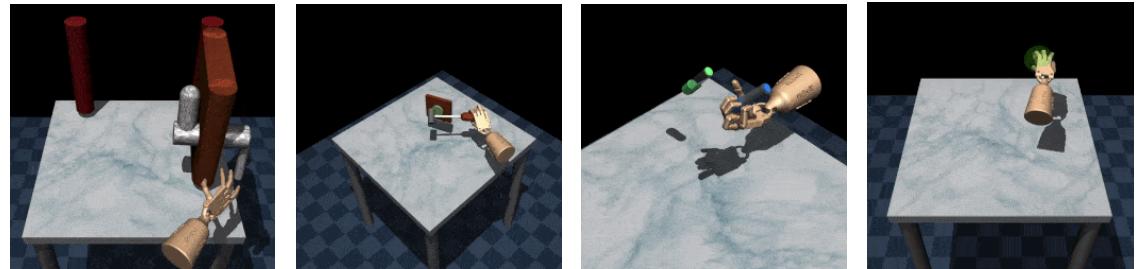
+ additional **mixed** datasets that are more difficult

mazes (2D maze & ant maze)



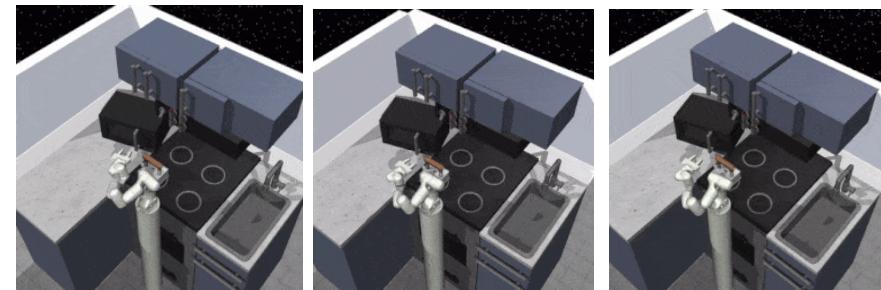
multi-task data, tests “stitching” ability

dexterous manipulation tasks



high dimensionality, human data, looks really awesome

whole-arm manipulation data



multi-task human data, realistic tasks

traffic control (Flow)

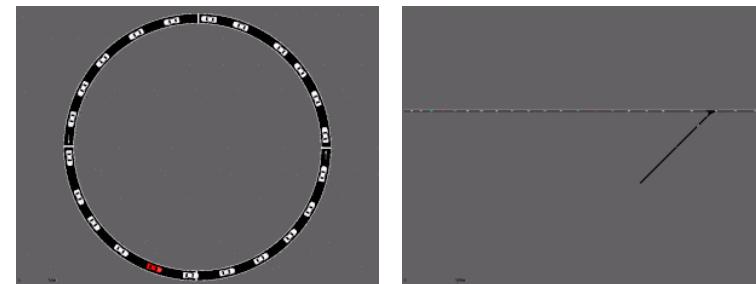
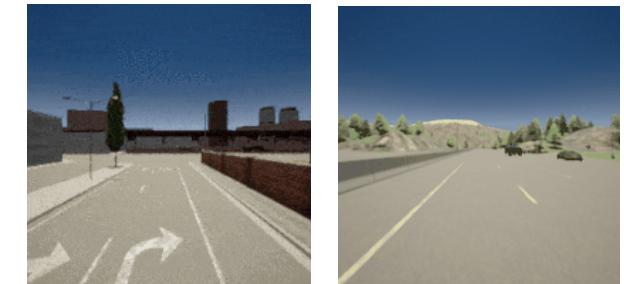


image-based driving (Carla)



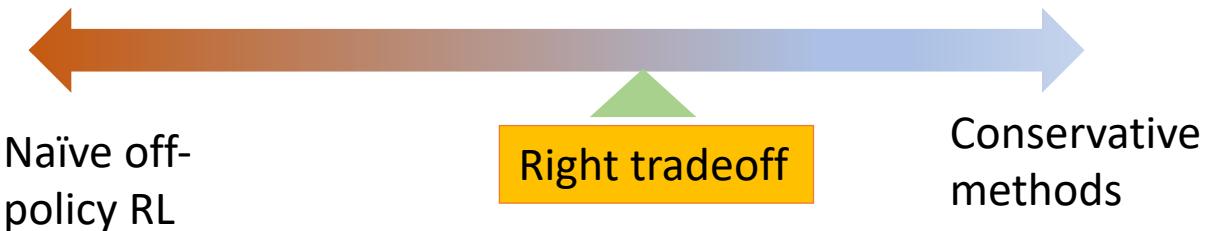
Other properties that are important

- **Stochasticity:** Common in real-world settings (economics, healthcare, education), often absent in commonly used benchmarks, high potential for impact from offline RL
- **Non-stationarity:** Things change over time, very common feature in the real world
- **“Risky” dataset biases:** If your driving data never shows a car crash, can you understand to avoid car crashes?
 - Some methods can handle this! But important to test
- **Partial observability:** Not discussed in this tutorial, but very important in the real world

1. Reinforcement learning primer
2. The offline reinforcement learning problem
3. Classic batch reinforcement learning algorithms
 - i. Importance-sampled policy gradients
 - ii. Least-squares temporal difference (LSTD)
4. Modern offline reinforcement learning algorithms
 - i. Policy constraint methods
 - ii. Model-based offline RL
 - iii. Value function regularization
 - iv. Uncertainty-based methods
 - v. Comparisons and discussion
5. Off-policy evaluation and model selection
 - i. Problem statement
 - ii. Direct methods and importance sampling methods
 - iii. Doubly robust methods
6. How should we evaluate offline RL methods?
 - i. When and why does offline RL work?
 - ii. How do we benchmark offline RL methods?
7. Conclusions, discussion, and open problems

Open Problems: Over-Conservatism?

- Methods covered so far act over-conservatively: there is no need to for conservatism if our learned models can generalize.
- Uncertainty-aware methods can be used to tradeoff conservatism for generalization



- But we need calibrated uncertainty estimates, calibrated across time!

Ways to measure and optimize uncertainty (consistent with RL)

Can we ever avoid conservatism?

“Expected case” scenarios

Open Problems: Hyperparameter Tuning

- **General Purpose OPE:** Do we need general OPE for hyperparameter tuning in offline RL?
Can we instead utilize a “validation error” metric during training?
Can we evaluate a subset of specific policies in OPE?
- **Goal:** Often we care about ranking policies, so do we actually need to train the policy return or value function?

Ranking correlation

Regret @ k

$$\text{Regret @ 1} = \max_{\pi \in \Pi} J(\pi) - J\left(\arg \max_{\pi \in \Pi} \hat{J}(\pi)\right)$$

Optimal policy in a set of policies

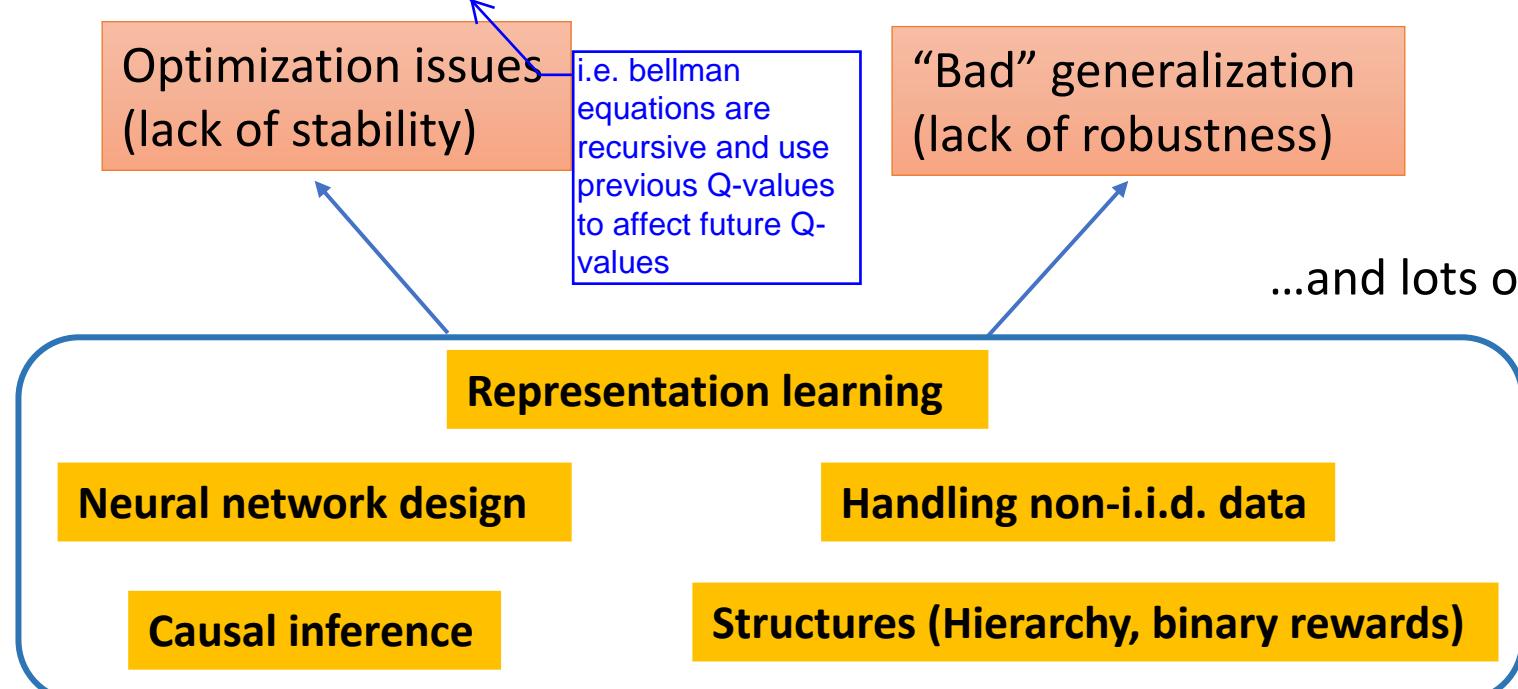
Pick out top k policies according to predicted return

Open Problems: Function Approximation

We want to leverage "good generalization" properties of deep neural networks in offline RL

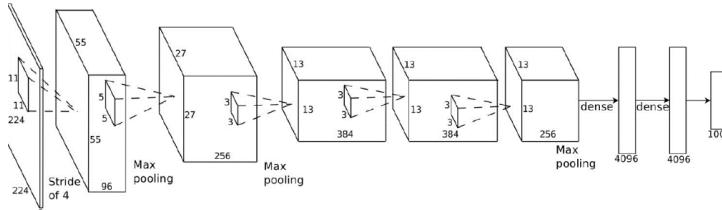
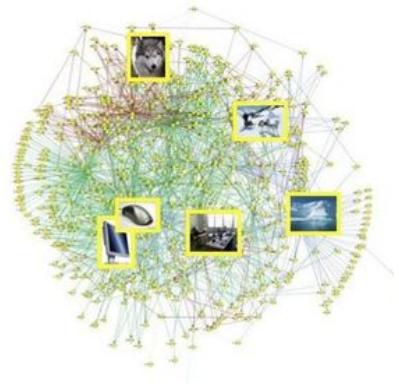
Function approximation can have drastic impact on the performance of RL in offline settings in particular

Using one to train itself can lead to "bad stuff"



...and lots of more open questions

Conclusion



enormous gulf



Offline RL methods aim to bridge this gulf by converting datasets to effective decision-making engines

Thank You!

Resources:

- Tutorial website (slides/video/more content):
<https://sites.google.com/view/offline-rl-tutorial>

- Offline RL survey paper:
<https://arxiv.org/abs/2005.01643>

- Hands-on coding assignment at:
<https://sites.google.com/view/offline-rl-tutorial>