# Offline Deep Reinforcement Learning Algorithms
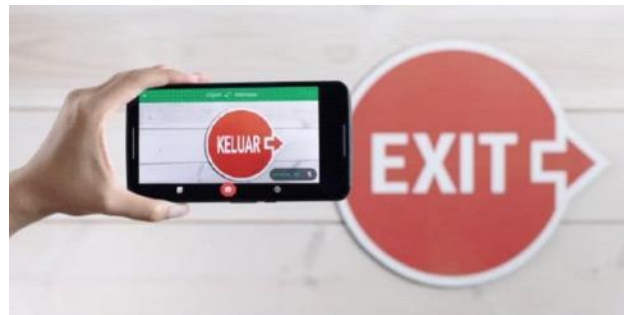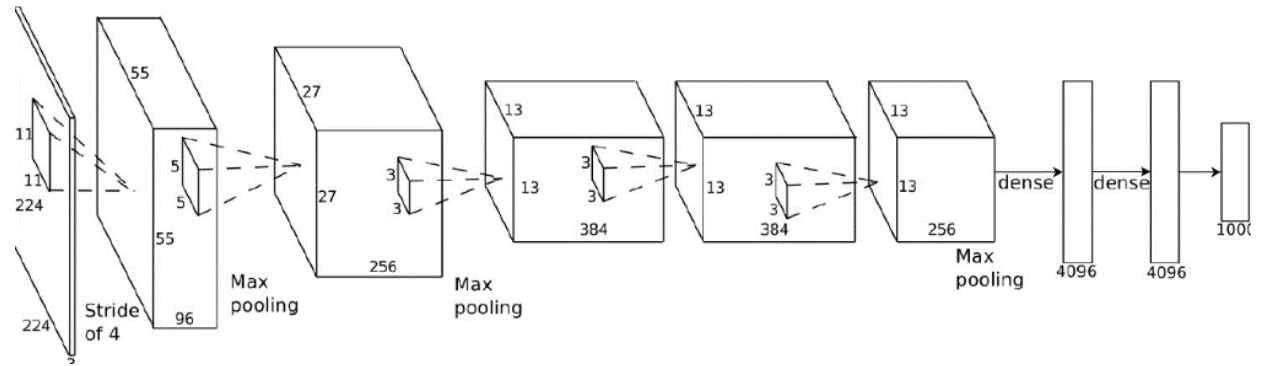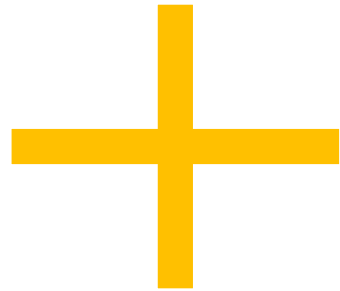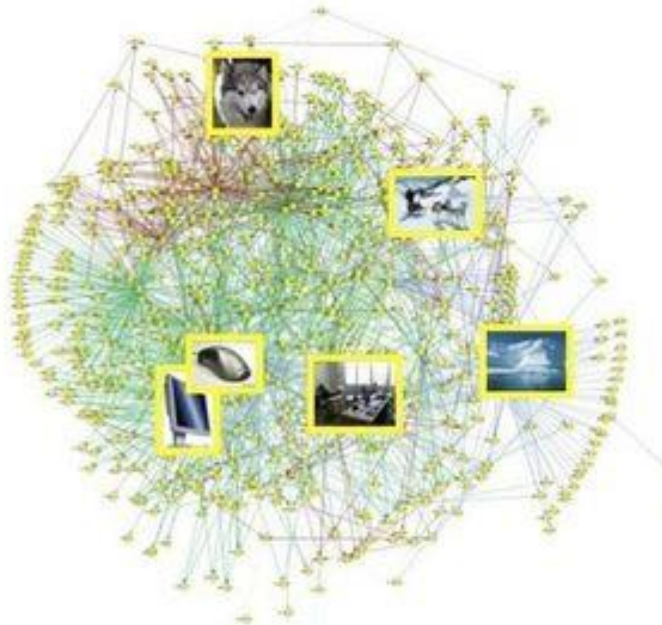
**Sergey Levine**

**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

this is done **many** times

enormous gulf

# Can we develop **data-driven** RL methods?



on-policy RL

off-policy RL

offline reinforcement learning

big datasets from past interaction

occasionally get more data

train for **many** epochs

Levine, Kumar, Tucker, Fu. **Offline Reinforcement Learning: Tutorial, Review, and Perspectives on Open Problems.** '20

Why is offline RL difficult?



How do we design offline RL algorithms?



Conservative Q-Learning

# Why is offline RL difficult?

# How do we design offline RL algorithms?

# Conservative Q-Learning

# Off-policy RL: a quick primer



RL objective: $\max_{\pi} \sum_{t=1}^{T} E_{\mathbf{s}_t, \mathbf{a}_t \sim \pi}[r(\mathbf{s}_t, \mathbf{a}_t)]$

Q-function: $Q^{\pi}(\mathbf{s}_t, \mathbf{a}_t) = \sum_{t'=t}^{T} E_{\mathbf{s}_{t'}, \mathbf{a}_{t'} \sim \pi}[r(\mathbf{s}_{t'}, \mathbf{a}_{t'}) | \mathbf{s}_t, \mathbf{a}_t]$

$\pi(\mathbf{a}|\mathbf{s}) = 1$ if $\mathbf{a} = \arg\max_{\mathbf{a}} Q^{\pi}(\mathbf{s}, \mathbf{a})$

$Q^{\star}(\mathbf{s}, \mathbf{a}) = r(\mathbf{s}, \mathbf{a}) + \max_{\mathbf{a}'} Q^{\star}(\mathbf{s}', \mathbf{a}')$

enforce this equation at all states!

minimize $\sum_{i}(Q(\mathbf{s}_i, \mathbf{a}_i) - [r(\mathbf{s}_i, \mathbf{a}_i) + \max_{\mathbf{a}_i'} Q(\mathbf{s}_i', \mathbf{a}_i')])^2$

minimize $\sum_{i}(Q(\mathbf{s}_i, \mathbf{a}_i) - y_i)^2$

This talk focuses entirely on **approximate dynamic programming** methods, but there are other methods too!

# Off-policy RL: a quick primer

$$Q(\mathbf{s}, \mathbf{a}) \leftarrow r(\mathbf{s}, \mathbf{a}) + \max_{\mathbf{a}'} Q(\mathbf{s}', \mathbf{a}')$$ ⟵———— don't need on-policy data for this!

off-policy Q-learning:

1. collect dataset $\{(\mathbf{s}_i, \mathbf{a}_i, \mathbf{s}'_i, r_i)\}$ using some policy, add it to $\mathcal{B}$

$K\times$

2. sample a batch $(\mathbf{s}_i, \mathbf{a}_i, \mathbf{s}'_i, r_i)$ from $\mathcal{B}$

3. minimize $\sum_i (Q(\mathbf{s}_i, \mathbf{a}_i) - [r(\mathbf{s}_i, \mathbf{a}_i) + \max_{\mathbf{a}'_i} Q(\mathbf{s}'_i, \mathbf{a}'_i)])^2$

$(\mathbf{s}, \mathbf{a}, \mathbf{s}', r)$

dataset of transitions
("replay buffer")

off-policy
Q-learning

$\pi(\mathbf{a}|\mathbf{s})$ (with exploration)

See, e.g.
Riedmiller, Neural Fitted Q-Iteration '05
Ernst et al., Tree-Based Batch Mode RL '05

# Does it work?



live data collection

Kalashnikov, Irpan, Pastor, Ibarz, Herzong, Jang, Quillen, Holly, Kalakrishnan, Vanhoucke, Levine. **QT-Opt: Scalable Deep Reinforcement Learning of Vision-Based Robotic Manipulation Skills**

# Does it work?



2x

4x speed



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

Kalashnikov, Irpan, Pastor, Ibarz, Herzong, Jang, Quillen, Holly, Kalakrishnan, Vanhoucke, Levine. **QT-Opt: Scalable Deep Reinforcement Learning of Vision-Based Robotic Manipulation Skills**

# What's the problem?

log scale (massive overestimation)

amount of data

**Hypothesis 1:** Overfitting



how well it does

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

**Hypothesis 2:** Training data is not good

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

Kumar, Fu, Tucker, Levine. **Stabilizing Off-Policy Q-Learning via Bootstrapping Error Reduction.** NeurIPS '19

Aviral Kumar

Justin Fu

# Distribution shift in a nutshell

Example empirical risk minimization (ERM) problem:

usually we are not worried – neural nets generalize well!

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

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

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

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

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

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

# Where do we suffer from distribution shift?

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

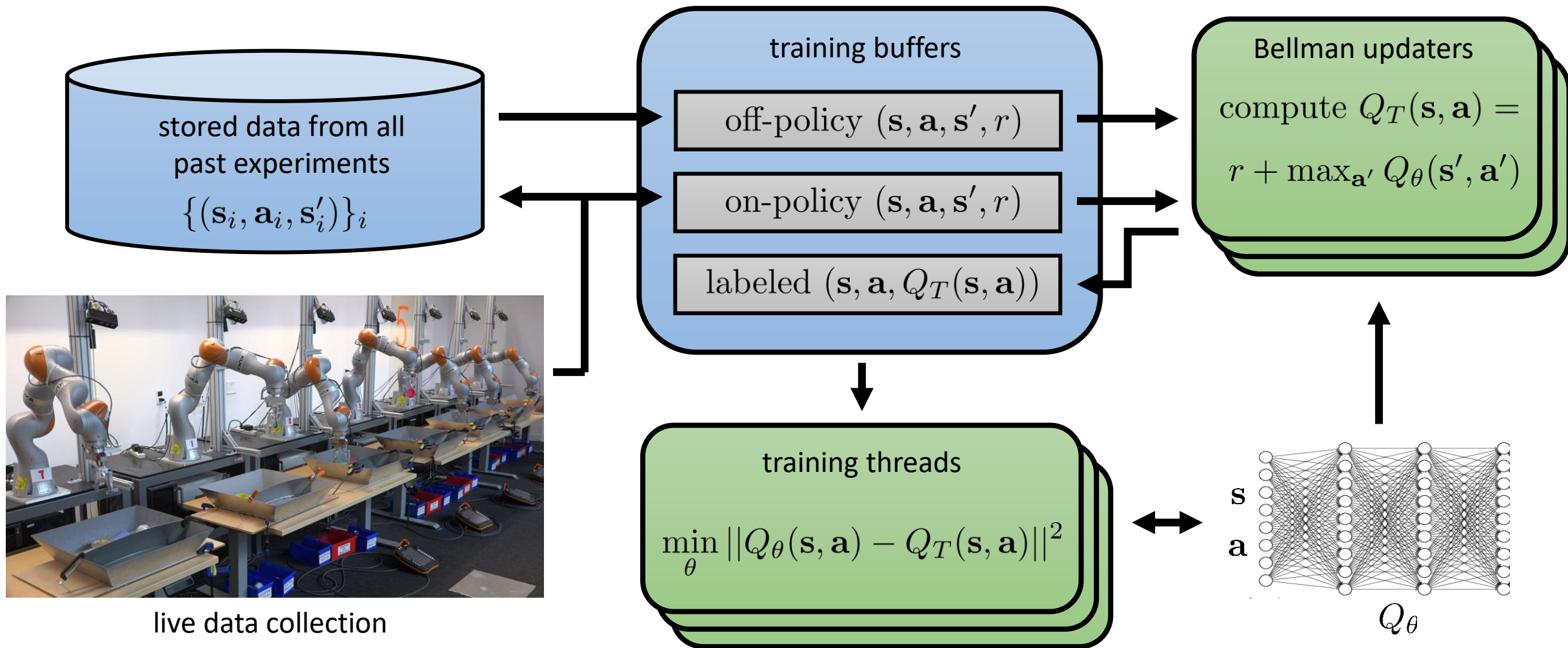$$Q(\mathbf{s}, \mathbf{a}) \leftarrow \underbrace{r(\mathbf{s}, \mathbf{a}) + E_{\mathbf{a}' \sim \pi_{\text{new}}}[Q(\mathbf{s}', \mathbf{a}')]}_{y(\mathbf{s}, \mathbf{a})}$$

what is the objective?

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

behavior policy

target value

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

how often does *that* happen?

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 $\mathbf{x}^{\star} \leftarrow \arg\max_{\mathbf{x}} f_{\theta}(\mathbf{x})$?)

how well it does

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

Kumar, Fu, Tucker, Levine. **Stabilizing Off-Policy Q-Learning via Bootstrapping Error Reduction.** NeurIPS '19

# Why is offline RL difficult?

# How do we design offline RL algorithms?

# Conservative Q-Learning

# How do prior methods address this?

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

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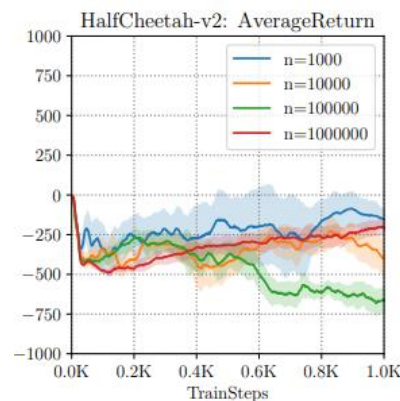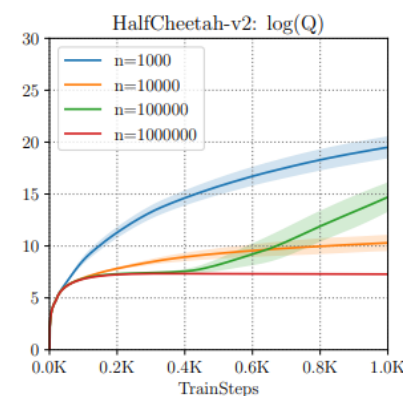
This solves distribution shift, right?

No more erroneous values?

can partially mitigate with **support** constraint (see Kumar et al. '19 "BEAR")

**Issue 1:** This might be **way** too conservative

**Issue 2:** Estimating the behavior policy is difficult

"policy constraint" method

**very** old idea (but it had no single name?)

Todorov et al. [passive dynamics in linearly-solvable MDPs]

Kappen et al. [KL-divergence control, etc.]

trust regions, covariant policy gradients, natural policy gradients, etc.

used in some form in recent papers:

Fox et al. '15 ("Taming the Noise…")

Fujimoto et al. '18 ("Off Policy…")

Jaques et al. '19 ("Way Off Policy…")

Kumar et al. '19 ("Stabilizing…")

Wu et al. '19 ("Behavior Regularized…")

Levine, Kumar, Tucker, Fu. **Offline Reinforcement Learning: Tutorial, Review, and Perspectives on Open Problems.** '20

# How bad is it?

**Issue 2:** Estimating the behavior policy is difficult

**Experiment:** online finetuning from offline initialization



see also:
Ghasemipour et al., EMaQ: Expected-Max Q-Learning Operator for Simple Yet Effective Offline and Online RL, '20

➢ More **powerful behavior policy** models lead to **improvement**, implying behavior policy modeling is a **major bottleneck**



offline training    online training

online training

Nair, Dalal, Gupta, Levine. **Accelerating Online Reinforcement Learning with Offline Datasets.** '20

# Avoiding behavior policies with **implicit** constraints

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

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

straightforward to
show via duality

**See also:**
Peters et al. (REPS)
Rawlik et al. ("psi-learning")
...many follow-ups

approximate via **weighted** max likelihood!

$$\pi_{\text{new}}(\mathbf{a}|\mathbf{s}) = \arg\max_\pi E_{(\mathbf{s},\mathbf{a})\sim\pi_\beta}\left[\log\pi(\mathbf{a}|\mathbf{s})\frac{1}{Z(\mathbf{s})}\exp\left(\frac{1}{\lambda}A^{\pi_{\text{old}}}(\mathbf{s},\mathbf{a})\right)\right]$$

samples from dataset
$\mathbf{a}\sim\pi_\beta(\mathbf{a}|\mathbf{s})$

critic can be used
to give us this

but maybe we can solve the overestimation problem at the **root**?

Peng*, Kumar*, Levine. **Advantage-Weighted Regression.** '19

Nair, Dalal, Gupta, Levine. **Accelerating Online Reinforcement Learning with Offline Datasets.** '20

# Why is offline RL difficult?

# How do we design offline RL algorithms?

# Conservative Q-Learning

# What about those Q-value errors?



how well it does

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

$$\hat{Q}^{\pi} = \arg \min_{Q} \max_{\mu} \alpha E_{\mathbf{s}\sim D, \mathbf{a}\sim\mu(\mathbf{a}|\mathbf{s})}[Q(\mathbf{s}, \mathbf{a})] \Big\}\text{— term to push down big Q-values}$$

regular objective $\Big\{ +E_{(\mathbf{s},\mathbf{a},\mathbf{s}')\sim D}\left[(Q(\mathbf{s}, \mathbf{a}) - (r(\mathbf{s}, \mathbf{a}) + E_{\pi}[Q(\mathbf{s}', \mathbf{a}')]))^2\right]$

can show that $\hat{Q}^{\pi} \leq Q^{\pi}$ for large enough $\alpha$

true Q-function

# Learning with Q-function lower bounds

Algorithm:

1. Learn $\hat{Q}^\pi$ for current $\pi$  such that $\hat{Q}^\pi \leq Q^\pi$

2. $\pi \leftarrow \arg\max_{\pi_{\text{new}}} E_{\pi_{\text{new}}}[\hat{Q}^\pi]$

A *better* bound:

always pushes Q-values down          push up on (**s**, **a**) samples in data

$$\hat{Q}^\pi = \arg\min_Q \max_\mu \alpha E_{\mathbf{s} \sim D, \mathbf{a} \sim \mu(\mathbf{a}|\mathbf{s})}[Q(\mathbf{s},\mathbf{a})] - \alpha E_{(\mathbf{s},\mathbf{a}) \sim D}[Q(\mathbf{s},\mathbf{a})]$$

$$+ E_{(\mathbf{s},\mathbf{a},\mathbf{s}') \sim D}\left[\left(Q(\mathbf{s},\mathbf{a}) - (r(\mathbf{s},\mathbf{a}) + E_\pi[Q(\mathbf{s}',\mathbf{a}')])\right)^2\right]$$

no longer guaranteed that $\hat{Q}^\pi(\mathbf{s},\mathbf{a}) \leq Q^\pi(\mathbf{s},\mathbf{a})$ *for all* $(\mathbf{s},\mathbf{a})$

but guaranteed that $E_{\pi(\mathbf{a}|\mathbf{s})}[\hat{Q}^\pi(\mathbf{s},\mathbf{a})] \leq E_{\pi(\mathbf{a}|\mathbf{s})}[Q^\pi(\mathbf{s},\mathbf{a})]$ *for all* $\mathbf{s} \in D$

Kumar, Zhou, Tucker, Levine. **Conservative Q-Learning for Offline Reinforcement Learning.** '20

Aviral
Kumar

# The conservative Q-learning (CQL) bound

minimize the **big** Q-values

maximize Q-values of state-action pairs in data

$$\hat{Q}^\pi_{\text{CQL}} = \arg\min_Q \max_\mu \alpha E_{\mathbf{s}\sim D, \mathbf{a}\sim\mu(\mathbf{a}|\mathbf{s})}[Q(\mathbf{s},\mathbf{a})] - \alpha E_{(\mathbf{s},\mathbf{a})\sim D}[Q(\mathbf{s},\mathbf{a})]$$

$$+\frac{1}{2}E_{(\mathbf{s},\mathbf{a},\mathbf{s}')\sim D}\left[(Q(\mathbf{s},\mathbf{a}) - (r(\mathbf{s},\mathbf{a}) + E_\pi[Q(\mathbf{s}',\mathbf{a}')]))^2\right]$$

**Theorem 3.2** (Equation 2 results in a tighter lower bound). *The value of the policy under the Q-function from Equation 2, $\hat{V}^\pi(\mathbf{s}) = \mathbb{E}_{\pi(\mathbf{a}|\mathbf{s})}[\hat{Q}^\pi(\mathbf{s},\mathbf{a})]$, lower-bounds the true value of the policy obtained via exact policy evaluation, $V^\pi(\mathbf{s}) = \mathbb{E}_{\pi(\mathbf{a}|\mathbf{s})}[Q^\pi(\mathbf{s},\mathbf{a})]$, when $\mu = \pi$, according to:*

concentration constant

$$\forall \mathbf{s}, \quad \hat{V}^\pi(\mathbf{s}) \le V^\pi(\mathbf{s}) - \alpha\left(I - \gamma P^\pi\right)^{-1}\mathbb{E}_{\pi(\mathbf{a}|\mathbf{s})}\left[\frac{\pi(\mathbf{a}|\mathbf{s})}{\hat{\pi}_\beta(\mathbf{a}|\mathbf{s})} - 1\right](\mathbf{s}) + \left(I - \gamma P^\pi\right)^{-1}\frac{C_{r,T,\delta}R_{\max}}{(1-\gamma)}.$$

pessimism due to regularizer

accounts for sampling error

Kumar, Zhou, Tucker, Levine. **Conservative Q-Learning for Offline Reinforcement Learning.** '20

# Does the bound hold in practice?

Underestimation vs. overestimation

$$E[\hat{Q}(\mathbf{s}, \mathbf{a})] - E[Q(\mathbf{s}, \mathbf{a})]$$

from Monte Carlo estimation

| 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 |

all prior methods have positive errors = wild optimism

CQL **always** has negative errors = pessimism

Kumar, Zhou, Tucker, Levine. **Conservative Q-Learning for Offline Reinforcement Learning.** '20

# 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



simulation & human data from Rajeswaran et al.

**Stitching:** data where dynamic programming can find much better solutions

**Realistic tasks**

Fu, Kumar, Nachum Tucker, Levine. **D4RL: Datasets for Data-Driven Deep Reinforcement Learning.** '20

Justin Fu

Aviral Kumar

# How does CQL compare?

| Task Name | QR-DQN | REM | CQL($\mathcal{H}$) |
|---|---|---|---|
| Pong (1%) | -13.8 | -6.9 | **19.3** |
| Breakout | 7.9 | 11.0 | **61.1** |
| Q*bert | 383.6 | | **14012.0** |
| Seaquest | 672.9 | 499.8 | **79.4** |
| Asterix* | 166.3 | 386.5 | **592.4** |

**1.5 – 6x** *better*

baseline: just clone the data

nothing works on the harder mazes?

nothing beats behavioral cloning?

| Domain | Task Name | BC | SAC | BEAR | BRAC-p | BRAC-v | CQL($\mathcal{H}$) | CQL($\rho$) |
|---|---|---|---|---|---|---|---|---|
| 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- | | | | | | **53.7** | 5.1 |
| | antmaze-large- | | | | 0.0 | | **15.8** | 3.2 |
| | antmaze-large- | | | | 0.0 | 0.0 | **14.9** | 2.3 |
| Adroit | pen-human | 34.4 | 6.3 | -1.0 | 8.1 | 0.6 | 37.5 | **55.8** |
| | hammer-human | | | | 0.3 | | **4.4** | 2.1 |
| | door-human | | | | -0.3 | | **9.9** | 9.1 |
| | relocate-human | | | | -0.3 | -0.3 | 0.20 | **0.35** |
| | pen-cloned | | | | 1.6 | -2.5 | 39.2 | 40.3 |
| | hammer-cloned | 0.8 | 0.2 | 0.3 | 0.3 | | 2.1 | **5.7** |
| | door-cloned | | | | -0.1 | | 0.4 | **3.5** |
| | relocate-cloned | | | | -0.3 | -0.3 | -0.1 | **-0.1** |
| Kitchen | kitchen-complete | | | | 0.0 | 0.0 | **43.8** | 31.3 |
| | kitchen-partial | | | | 0.0 | | **49.8** | 50.1 |
| | kitchen-undirected | 47.5 | 2.5 | 47.2 | 0.0 | 0.0 | **51.0** | 52.4 |

*"infinitely" better*

**1.5-3x** *better*

**up to 5x** *better*

**1.1 – 1.3x** *better*

CQL seems to work pretty well on many tasks!

And we seem to know *why* it works!

But there is still plenty of room for improvement…

Kumar, Zhou, Tucker, Levine. **Conservative Q-Learning for Offline Reinforcement Learning.** '20

- Offline RL is quite difficult, but has **enormous promise**, and initial results suggest it can be **extremely powerful**

- Effective (dynamic programming) offline RL methods can be implemented by imposing **constraints** on the policy, perhaps implicitly

- Learning a lower bound Q-function (i.e., conservative Q-learning) can **substantially** improve offline RL performance

$$\hat{Q}^{\pi}_{\mathrm{CQL}} = \arg\min_{Q}\max_{\mu} \alpha E_{\mathbf{s}\sim D, \mathbf{a}\sim\mu(\mathbf{a}|\mathbf{s})}[Q(\mathbf{s},\mathbf{a})] - \alpha E_{(\mathbf{s},\mathbf{a})\sim D}[Q(\mathbf{s},\mathbf{a})]$$
$$+\frac{1}{2}E_{(\mathbf{s},\mathbf{a},\mathbf{s}')\sim D}\left[\left(Q(\mathbf{s},\mathbf{a})-(r(\mathbf{s},\mathbf{a})+E_{\pi}[Q(\mathbf{s}',\mathbf{a}')])\right)^2\right]$$

Kumar, Fu, Tucker, Levine. **Stabilizing Off-Policy Q-Learning via Bootstrapping Error Reduction.** NeurIPS '19

Nair, Dalal, Gupta, Levine. **Accelerating Online Reinforcement Learning with Offline Datasets.** '20

Kumar, Zhou, Tucker, Levine. **Conservative Q-Learning for Offline Reinforcement Learning.** '20

Fu, Kumar, Nachum Tucker, Levine. **D4RL: Datasets for Data-Driven Deep Reinforcement Learning.** '20