Rishabh Agarwal, Dale Schuurmans, Mohammad Norouzi

HOW I LEARNED TO STOP WORRYING AND LOVE OFFLINE RL

An Optimistic Perspective on Offline Reinforcement Learning







Reinforcement Learning with Online Interactions





Offline Reinforcement Learning





Offline RL: Then and Now

- (Then) Offline RL was previously known as "Batch RL".
 - Our paper was the first to use the term "offline RL" and popularize it!
- (Then) ICLR'20 review: "Why are offline algorithms necessary? I am not sure this paper is relevant to the community of standard RL."
 - Eventually, accepted at ICML'20.
- (Now) Organizing 1st offline RL workshop at NeurIPS 2020 offline-rl-neurips.github.io
 - 50+ accepted papers at the workshop

Offline RL on Atari 2600



Talk



What makes Deep Learning Successful?

Expressive function approximators





What makes Deep Learning Successful?

Expressive function approximators



Powerful learning algorithms



What makes Deep Learning Successful?

Expressive function approximators



Large and Diverse Datasets



Powerful learning algorithms





Expressive function approximators



Good learning algorithms e.g., actor-critic, approx DP

Expressive function approximators



Good learning algorithms e.g., actor-critic, approx DP



Google Research

Expressive function approximators



Good learning algorithms e.g., actor-critic, approx DP

Interactive Environments



Google Research





Robotics

Dasari, Ebert, Tian, Nair, Bucher, Schmeckpeper, .. Finn. RoboNet: Large-Scale Multi-Robot Learning.
Yu, Xian, Chen, Liu, Liao, Madhavan, Darrell. BDD100K: A Large-scale Diverse Driving Video Database.





Recommender Systems

Dasari, Ebert, Tian, Nair, Bucher, Schmeckpeper, .. Finn. RoboNet: Large-Scale Multi-Robot Learning.
Yu, Xian, Chen, Liu, Liao, Madhavan, Darrell. BDD100K: A Large-scale Diverse Driving Video Database.





Robotics





Autonomous Driving

Recommender Systems

Dasari, Ebert, Tian, Nair, Bucher, Schmeckpeper, ..., Finn. RoboNet: Large-Scale Multi-Robot Learning.
Yu, Xian, Chen, Liu, Liao, Madhavan, Darrell. BDD100K: A Large-scale Diverse Driving Video Database.

Vikas Atta



Google Research



Reinforcement Learning with Online Interactions





Offline Reinforcement Learning







Offline RL can help:

• **Pretrain** agents on existing logged data.





Offline Reinforcement Learning

Reinforcement Learning with Online Interactions







Offline RL can help:

- Pretrain agents on existing logged data.
- **Evaluate** RL algorithms on the basis of exploitation alone on common datasets.

Reinforcement Learning with Online Interactions





Offline Reinforcement Learning







Offline RL can help:

- Pretrain the agents on existing logged data.
- Evaluate RL algorithms on the basis of exploitation alone on common datasets.
- Deliver real-world **impact**.

Reinforcement Learning with Online Interactions





Offline Reinforcement Learning







But .. Offline RL is Challenging!



Distribution mismatch



But .. Offline RL is Challenging!



No New Corrective Feedback



But .. Offline RL is Challenging!

Fully Off-Policy



Standard RL fails in the Offline setting?

Off-Policy Deep Reinforcement Learning without Exploration

Scott Fujimoto¹² David Meger¹² Doina Precup¹²

Abstract

Many practical applications of reinforcement learning constrain agents to learn from a fixed batch of data which has already been gathered, without offering further possibility for data collection. In this paper, we demonstrate that due to

require further interactions with the environment to compensate (Hester et al., 2017; Sun et al., 2018; Cheng et al., 2018). On the other hand, batch reinforcement learning offers a mechanism for learning from a fixed dataset without restrictions on the quality of the data.

Most modern off-policy deep reinforcement learning al-

Ofir Nachum

Google Research

Behavior Regularized Offline Reinforcement Learning

Yifan Wu* Carnegie Mellon University yw4@cs.cmu.edu

George Tucker Google Research gjt@google.com ofirnachum@google.com

Abstract

In reinforcement learning (RL) research, it is common to assume access to direct online interactions with the environment. However in many real-world applications, access to the environment is limited to a fixed offline dataset of logged experience. In such settings, standard RL algorithms have been shown to diverge or otherwise yield poor performance. Accordingly, recent work has suggested a number of remedies to these issues. In this work, we introduce a general framework, behavior regularized actor critic (BRAC), to empirically evaluate recently proposed methods as well as a number of simple baselines across a variety of offline continuous control tasks. Surprisingly, we find that many of the technical complexities introduced in recent methods are unnecessary to achieve strong performance. Additional ablations provide insights into which design choices matter most in the offline RL setting.¹

KEEP DOING WHAT WORKED: BEHAVIOR MODELLING PRIORS FOR OFFLINE REIN-FORCEMENT LEARNING

Noah Y. Siegel, Jost Tobias Springenberg, Felix Berkenkamp, Abbas Abdolmaleki, Michael Neunert, Thomas Lampe, Roland Hafner, Nicolas Heess, Martin Riedmiller

> DeepMind {siegeln}@google.com

ABSTRACT

Off-policy reinforcement learning algorithms promise to be applicable in settings where only a fixed data-set (batch) of environment interactions is available and no new experience can be acquired. This property makes these algorithms appealing

Stabilizing Off-Policy Q-Learning via Bootstrapping **Error Reduction**

Aviral Kumar* UC Berkeley aviralk@berkeley.edu

gjt@google.com

George Tucker Google Brain

Justin Fu* UC Berkeley justinjfu@eecs.berkeley.edu

Sergey Levine UC Berkeley, Google Brain svlevine@eecs.berkeley.edu

Abstract

Off-policy reinforcement learning aims to leverage experience collected from prior policies for sample-efficient learning. However, in practice, commonly used off-policy approximate dynamic programming methods based on Q-learning and



Standard RL fails in the Offline setting ..



choices matter most in the offline RL setting.¹

off-policy approximate dynamic programming methods based on Q-learning and

Can standard off-policy RL succeed in the offline setting?



Offline RL on Atari 2600



Train 5 DQN (Nature) agents on 60 Atari games with sticky actions for 200 million frames.



Offline RL on Atari 2600



Save all (observation, action, next observation, reward) tuples encountered to **DQN Replay Dataset**. Total of 300 datasets, 5 per game.



Offline RL on Atari 2600



Train offline agents using DQN Replay Dataset without any further environment interactions.



Offline DQN on DQN Replay Dataset



Does Offline DQN work?





Let's try recent off-policy methods!



Distributional RL uses Z(s, a), a distribution over returns, instead of the *Q*-function.

$$Z(s,a;\theta) := \frac{1}{K} \sum_{i=1}^{K} \delta_{\theta_i(s,a)}$$
$$Q(s,a;\theta) := \mathbb{E}[Z] = \frac{1}{K} \sum_{i=1}^{K} \theta_i(s,a)$$

Does Offline QR-DQN work?



Does Offline DQN work?



Does Offline QR-DQN work?





Developing Robust Offline RL algorithms

- > Emphasis on Generalization
 - Given a fixed dataset, generalize to unseen states during evaluation.

Developing Robust Offline RL algorithms

- Emphasis on Generalization
 - Given a fixed dataset, generalize to unseen states during evaluation.
- > Ensemble of Q-estimates:
 - Ensembling, Dropout widely used for improving generalization.

Ensemble-DQN



Train multiple (linear) Q-estimates with different random initialization.

Google Research



Does Offline Ensemble-DQN work?



Offline DQN



Developing Robust Offline RL algorithms

- > Emphasis on Generalization
 - Given a fixed dataset, generalize to unseen states during evaluation.
- \succ Q-learning as constraint satisfaction:

•
$$\forall (s, a, s', r) : Q^*(s, a) = r + max_{a'} Q^*(s', a')$$



Random Ensemble Mixture (REM)



Minimize TD error on random (per minibatch) convex combination of multiple Q-estimates.

REM vs QR-DQN





QR-DQN

Google Research

Offline Stochastic Atari Results



Scores averaged over 5 runs of offline agents trained using DQN replay data across 60 Atari games for 5X gradient steps. Offline REM surpasses gains from online C51 and offline QR-DQN.

Google Research

Offline REM vs. Baselines



Does Online REM work?



Average normalized scores of online agents trained for 200 million game frames. Multi-network REM with 4 Q-functions performs comparably to QR-DQN.

Important Factors in Offline RL

An Optimistic Perspective on Offline Reinforcement Learning

Key Factor in Success: Offline Dataset Size



Randomly subsample N% of frames from 200 million frames for offline training.

Key Factor in Success: Offline Dataset Diversity



Subsample first 10% of total frames (20 million) for offline training -- much lower quality data.

Choice of Algorithm: Offline Continuous Control



Google Research

Offline agents trained using full experience replay of DDPG on MuJoCo environments.



Overfitting in Offline RL: Number of Gradient Updates



Average online scores of offline agents trained on 5 games using logged DQN replay data for 5X gradient steps compared to online DQN.

Overfitting Underfitting in Offline RL: Number of Gradient Updates



Reason: Implicit Regularization of gradient descent ("Preference for simpler solutions") is amplified by bootstrapping (learning a guess from guess) in RL.

Implicit Under-Parameterization in Deep RL



Implicit Under-Parameterization in Deep RL



Q-network implicitly behaves as low capacity network!

Offline RL for Robotics

Scaling data-driven robotics with reward sketching and batch reinforcement learning

Serkan Cabi, Sergio Gómez Colmenarejo, Alexander Novikov, Ksenia Konyushkova, Scott Reed, Rae Jeong, Konrad Żołna, Yusuf Aytar, David Budden, Mel Vecerik, Oleg Sushkov, David Barker, Jonathan Scholz, Misha Denil, Nando de Freitas, Ziyu Wang

Abstract—By harnessing a growing dataset of robot experience, we learn control policies for a diverse and increasing set of related manipulation tasks. To make this possible, we introduce reward sketching: an effective way of eliciting human preferences to learn the reward function for a new task. This reward function is then used to retrospectively annotate all historical data, collected for different tasks, with predicted rewards for the new task. The resulting massive annotated dataset can then be used to learn manipulation policies with batch reinforcement learning (RL) from visual input in a completely off-line way, i.e. without interaction with the real robot. This approach makes it possible to scale up RL in robotics, as we no longer need to run the robot for each step of learning. We show that the trained batch RL agents, when deployed in real robots, can perform a variety of challenging tasks involving multiple interactions among nigid on deformable objects Manager that display a significant



Future Work

"The potential for off-policy learning remains tantalizing, the best way to achieve it still a mystery." - Sutton & Barto

G. Barto

Reinforcement

Offline RL: Future Work

• Rigorous characterization of role of generalization in offline RL

- Rigorous characterization of role of generalization in offline RL
- Benchmarking with various data collection strategies
 - Subsampling **DQN Replay Dataset** (e.g., first / last *k* million frames)

- Rigorous characterization of role of generalization in offline RL
- Benchmarking with various data collection strategies
 - Subsampling DQN-replay datasets (e.g., first / last *k* million frames)
- Offline Evaluation / Hyperparameter Tuning

- Rigorous characterization of role of generalization in offline RL
- Benchmarking with various data collection strategies
 - Subsampling DQN-replay datasets (e.g., first / last *k* million frames)
- Offline Evaluation / Hyperparameter Tuning
- Self-supervised / Model-based RL approaches

- Rigorous characterization of role of generalization in offline RL
- Benchmarking with various data collection strategies
 - Subsampling DQN-replay datasets (e.g., first / last *k* million frames)
- Offline Evaluation / Hyperparameter Tuning
- Self-supervised / Model-based RL approaches
- Combining REM with behavior regularization (BCQ, SPIBB, CQL etc.)

TL;DR

- Standard RL algorithms (e.g. REM, QR-DQN), trained on sufficiently large and diverse datasets, perform quite well in the offline setting.
- Offline RL provides a standardized setup for:
 - Isolating *exploitation* from exploration
 - Developing *sample efficient* and *stable algorithms*
 - Pretrain RL agents on logged data

Thank you!

Code, dataset, blog and paper at offline-rl.qithub.io

An Optimistic Perspective on Offline Reinforcement Learning