Offline RL for Real-Robot Pre-Training and Fine-Tuning

Aviral Kumar



Offline Reinforcement Learning



No unsafe or costly exploration

Potential to bring generalization benefits of supervised learning

What this Picture Actually Looks Like



Learning from Diverse Robot Datasets

Toy kitchen 1



Toy Kitchen 3



Toy Sink 1







Put Pear In Bowl Put Sweet Potato In Pot

Toy Kitchen 4

Toy Sink 3

t Detergent into Drying





Pre-training on broad data





Pre-Training for Robots Using Offline RL



Ebert*, Yang* et al. Bridge data: Boosting Generalization of Robotic Skills. RSS 2022.

Ingredient 1: Conservative Q-Learning

The issue in offline RL is erroneous Q-values at out-of-distribution (OOD) actions



K., Zhou, Tucker, Levine. Conservative Q-Learning for Offline Reinforcement Learning. NeurIPS 2020



Ingredient 3: Rewards & Checkpoint Selection

Rewards

Sparse, binary rewards, +1 at the end of the trajectory

Important: the binary values matter (-1, +10)

Checkpoint selection

Worst case: impossible!

But can use the knowledge that data is "expert"



Summary of Ingredients in PTR

Ingredient 1: An Offline RL Algorithm (CQL)

Ingredient 2: A high-capacity architecture (ResNet + group normalization + action duplication + learned spatial embeddings)

Ingredient 3: Reward functions + checkpoint selection heuristic

Now some Empirical Results....

Task: Solving A Task in A New Domain





2. Fine-Tune on Target Domain Data:1 door, 10 demonstrations

1. Pre-Train on Bridge Data, 12 doors 800 demonstrations

Results: Solving A Task in A New Domain

Method: Imitation (Best prior method)



Method: PTR (Ours)



Task: Solving New Tasks in New Domains





10 target

demonstrations







Results: Solving New Tasks in New Domains

















Some Quantitative Results

Joint training vs pre-training??

		BC finetuning			Joint training Tar		Target	data only	Pre-train. rep. + BC finetune	
Task	PTR (Ours)	BC (fine.)	Autoreg. BC	BeT	COG	BC	CQL	BC	R3M	MAE
Take croissant from metal bowl Put sweet potato on plate Place knife in pot Put cucumber in pot	7/10 7/20 4/10 5/10	3/10 1/20 2/10 0/10	5/10 1/20 2/10 1/10	1/10 0/20 0/10 0/10	4/10 0/20 1/10 2/10	4/10 0/20 3/10 1/10	0/10 0/20 3/10 0/10	1/10 0/20 0/10 0/10	1/10 0/20 0/10 0/10	3/10 1/20 0/10 0/10
		Imitation (using transformers, auto-regressive)						Self-supervised pre-training from videos / bridge data		
					Ļ					
		Bette			Representation learning					

Takeaway: Offline RL learns useful representations + better fine-tuning

Scaling Curve And Analysis



Better performance with larger networks!

Why would RL enable better performance...

....when the data is collected via human teleoperation?

locate the pot accurately

PTR Success: Places Cucumber in Pot

Preview: Value-functions can learn what's critical!

Qualitative Comparison of BC (finetune) and PTR

Task: Take Croissant from Metal Bowl





Analysis of Why PTR Outperforms Imitation



Test: Run weighted BC, where weights come from the learned Q-function!

Task	BC (finetune)	TR (Ours)	Advantage-weighted BC (finetune)
Put cucumber in pot	0/10	5/10	5/10
Take croissant from metal bowl	3/10	7/10	6/10

K., Hong, Singh, Levine. When Should Offline RL Be Preferred Over BC? ICLR 2022.

Takeaways and Future Directions

Offline RL can be good for both representation learning and control, even with human demonstration data

Future Directions:

Extend to use videos and multi-robot data on more dexterous tasks

Goal specification: language? goals? reward learning?

> Workflows: "How should a practitioner tune this approach on their problem"

Thank You!

Paper: <u>https://arxiv.org/abs/2210.05178</u> Code: <u>https://github.com/Asap7772/PTR</u>

