Estimating Q(s,s') with Deep Deterministic Dynamics Gradients

Ashley D. Edwards, Himanshu Sahni, Rosanne Liu, Jane Hung, Ankit Jain, Rui Wang, Adrien Ecoffet, Thomas Miconi, Charles Isbell, Jason Yosinski







- Naturally learns well in redundant action spaces
- Transfers across permuted action spaces
- Enables learning from observations obtained from completely random policies

$$Q(s,s') \doteq r + \gamma \max_{s'' \in N(s')} Q(s',s'')$$

$$Q(s,s') \doteq r + \gamma \max_{s'' \in N(s')} Q(s',s'')$$

$$Q(s,s') = Q(s,s') + \alpha [r + \gamma \max_{s'' \in N(s)} Q(s',s'') - Q(s,s')]$$

$$Q(s, s') \doteq r + \gamma \max_{s'' \in N(s')} Q(s', s'')$$
$$Q(s, s') = Q(s, s') + \alpha [r + \gamma \max_{s'' \in N(s)} Q(s', s'') - Q(s, s')]$$
$$\tau(s) = \arg\max_{s' \in N(s)} Q(s, s')$$



- 1) Find neighboring states N(s)
- 2) Select state with largest value
- 3) Determine action



- 1) Find neighboring states N(s)
- 2) Select state with largest value
- 3) Determine action



- 1) Find neighboring states N(s)
- 2) Select state with largest value
- 3) Determine action

$$\tau(s) = \underset{s' \in N(s)}{\arg \max} Q(s, s')$$



<u>Control</u>

- 1) Find neighboring states N(s)
- 2) Select state with largest value
- 3) Determine action

$$\pi(s) = I(s, \tau(s))$$

- Q(s,s') stored in table
- Given N(s) and I(s,s')
- 11x11 gridworld
- 4 actions

		_	_	_		
						G



Example of equivalence of QSA and QSS

- 4 actions with probability of slipping
- Given N(s) and I(s,s')





Example of QSS with stochastic actions



(a) 25% (b) 50% (c) 75%

Example of QSS with stochastic actions



Example of QSS with stochastic actions

- 4 base actions with redundancy
- Learn I(s,s')





(a) QSA







Episode

30000

40000

50000

20000

40

0

-100

Average Episodic Reward

-400

-500

0

10000

400

QSS with redundant actions



- 11x11 gridworld
- 4 actions 0, 1, 2, 3
- Given N(s)
- Learn I(s,s')

- 11x11 gridworld
- 4 actions 3, 2, 1, 0
- Given N(s)
- Learn I(s,s')



QSS with permuted actions

- Experiments demonstrated properties of QSS
- Most problems cannot be solved in tabular setting
- How to learn in settings with large/continuous state spaces?

$$Q(s,s') \doteq r + \gamma \max_{s'' \in N(s')} Q(s',s'')$$

$$Q(s,s') = Q(s,s') + \alpha [r + \gamma \max_{s'' \in N(s)} Q(s',s'') - Q(s,s')]$$

$$\pi(s) = I(s,\tau(s))$$

$$\tau(s) = \underset{s' \in N(s)}{\operatorname{arg\,max}} Q(s,s')$$

$$Q(s,s') \doteq r + \gamma \max_{s'' \in N(s')} Q(s',s'')$$

$$Q(s,s') = Q(s,s') + \alpha[r + \gamma \max_{s'' \in N(s)} Q(s',s'') - Q(s,s')]$$

$$\pi(s) = I(s,\tau(s))$$

$$\tau(s) = \arg\max_{s' \in N(s)} Q(s,s')$$

$$Q(s,s') \doteq r + \gamma \max_{s'' \in N(s')} Q(s',s'')$$
$$Q(s,s') = Q(s,s') + \alpha [r + \gamma \max_{s'' \in N(s)} Q(s',s'') - Q(s,s')]$$
$$\pi(s) = I(s,\tau(s))$$
$$\tau(s) = \arg\max_{s' \in N(s)} Q(s,s')$$

$$Q(s, s') \doteq r + \gamma Q(s, \tau(s))$$

$$Q(s, s') = Q(s, s') + \alpha [r + \gamma \max_{s'' \in N(s)} Q(s', s'') - Q(s, s')]$$

$$\pi(s) = I(s, \tau(s))$$

$$\tau(s) = \underset{s' \in N(s)}{\operatorname{arg\,max}} Q(s, s')$$

$$Q(s, s') \doteq r + \gamma Q(s, \tau(s))$$

$$Q(s, s') = Q(s, s') + \alpha [r + \gamma \max_{s'' \in N(s)} Q(s', s'') - Q(s, s')]$$

$$\pi(s) = I(s, \tau(s))$$

$$\tau(s) = \underset{s' \in N(s)}{\operatorname{arg\,max}} Q(s, s')$$

$$Q(s, s') \doteq r + \gamma Q(s, \tau(s))$$

$$Q(s, s') = Q(s, s') + \alpha [r + \gamma Q(s', \tau(s')) - Q(s, s')]$$

$$\pi(s) = I(s, \tau(s))$$

$$\tau(s) = \underset{s' \in N(s)}{\arg \max} Q(s, s')$$

$$Q(s, s') \doteq r + \gamma Q(s, \tau(s))$$

$$Q(s, s') = Q(s, s') + \alpha [r + \gamma Q(s', \tau(s')) - Q(s, s')]$$

$$\pi(s) = I(s, \tau(s))$$

$$\tau(s) = \underset{s' \in N(s)}{\arg \max} Q(s, s')$$

$$Q(s, s') \doteq r + \gamma Q(s, \tau(s))$$

$$\mathcal{L}_{\theta} = \sum_{i} \|y - Q_{\theta_{i}}(s, s')\|$$
$$y = r + \gamma \min_{i=1,2} Q_{\theta_{i}'}(s', \tau_{\psi'}(s'))]$$

$$\pi(s) = I(s, \tau(s))$$

$$\tau(s) = \underset{s' \in N(s)}{\operatorname{arg\,max}} Q(s, s')$$

$$Q(s, s') \doteq r + \gamma Q(s, \tau(s))$$
$$\mathcal{L}_{\theta} = \sum_{i} \|y - Q_{\theta_{i}}(s, s')\|$$
$$y = r + \gamma \min_{i=1,2} Q_{\theta_{i}'}(s', \tau_{\psi'}(s'))]$$
$$\pi(s) = I(s, \tau(s))$$

$$\tau(s) = \underset{s' \in N(s)}{\operatorname{arg\,max}} Q(s, s')$$

$$Q(s, s') \doteq r + \gamma Q(s, \tau(s))$$
$$\mathcal{L}_{\theta} = \sum_{i} \|y - Q_{\theta_{i}}(s, s')\|$$
$$y = r + \gamma \min_{i=1,2} Q_{\theta_{i}'}(s', \tau_{\psi'}(s'))]$$
$$\mathcal{L}_{\omega} = \|I_{\omega}(s, s') - a\|$$
$$\tau(s) = \underset{s' \in N(s)}{\operatorname{arg\,max}} Q(s, s')$$

$$Q(s, s') \doteq r + \gamma Q(s, \tau(s))$$
$$\mathcal{L}_{\theta} = \sum_{i} \|y - Q_{\theta_{i}}(s, s')\|$$
$$y = r + \gamma \min_{i=1,2} Q_{\theta'_{i}}(s', \tau_{\psi'}(s'))]$$
$$\mathcal{L}_{\omega} = \|I_{\omega}(s, s') - a\|$$

$$\tau(s) = \underset{s' \in N(s)}{\operatorname{arg\,max}} Q(s, s')$$

$$Q(s, s') \doteq r + \gamma Q(s, \tau(s))$$
$$\mathcal{L}_{\theta} = \sum_{i} \|y - Q_{\theta_{i}}(s, s')\|$$
$$y = r + \gamma \min_{i=1,2} Q_{\theta'_{i}}(s', \tau_{\psi'}(s'))]$$
$$\mathcal{L}_{\omega} = \|I_{\omega}(s, s') - a\|$$

$$\mathcal{L}_{\psi} = -Q_{\theta}(s, \tau_{\psi}(s))$$

Continuous QSS via D3G



- 1) Find neighboring states N(s)
- 2) Select state with largest value
- 3) Determine action

$$s'_{\tau} \leftarrow \tau(s)$$

Continuous QSS via D3G



<u>Control</u>

- 1) Find neighboring states N(s)
- 2) Select state with largest value
- 3) Determine action

$$s'_{\tau} \leftarrow \tau(s)$$

Continuous QSS via D3G



- 1) Find neighboring states N(s)
- 2) Select state with largest value
- 3) Determine action

$$\pi(s) \leftarrow I(s, s_{\tau}')$$

- 11x11 gridworld
- Mujoco tasks
- I(s,s') learned

				_	_	
-						
						G





 $\mathcal{L}_{\psi} = -Q_{\theta}(s, \tau_{\psi}(s))$

critic_loss



 $\mathcal{L}_{\theta} = \sum \|y - Q_{\theta_i}(s, s')\|$ $y = r + \gamma \min_{i=1,2} Q_{\theta'_i}(s', \tau_{\psi'}(s'))]$



Model $au_{\psi}(s) = s_{\tau}'$

Inverse dynamics $I_{\omega}(s, s_{\tau}') = a$

Forward dynamics $f_{\phi}(s, a) = s'_f$



Model
$$au_{\psi}(s) = s'_{\tau}$$

Inverse dynamics $I_{\omega}(s, s'_{\tau}) = a$
Forward dynamics $f_{\phi}(s, a) = s'_{f}$

$$\mathcal{L}_{\psi} = -Q_{\theta}(s, C(s, \tau_{\psi}(s)) + \beta \| \tau_{\psi}(s) - C(s, \tau_{\psi}(s)) \|$$
$$y = r + \gamma \min_{i=1,2} Q_{\theta'_{i}}(s', C(s', \tau_{\psi'}(s')))$$



- 1) Select state with largest value
- 2) Determine action
- 3) Predict next state

$$s'_{\tau} \leftarrow \tau(s)$$



- 1) Select state with largest value
- 2) Determine action
- 3) Predict next state



 $s'_{\tau} \leftarrow \tau(s)$



- 1) Select state with largest value
- 2) Determine action
- 3) Predict next state



 $a \leftarrow I(s, s'_{\tau})$



- 1) Select state with largest value
- 2) Determine action
- 3) Predict next state



 $s'_f \leftarrow f(s, a)$



- 1) Select state with largest value
- 2) Determine action
- 3) Predict next state

$$s'_{\tau} \leftarrow \tau(s)$$



- 1) Select state with largest value
- 2) Determine action
- 3) Predict next state

$$s'_{\tau} \leftarrow \tau(s)$$













 $Q(s, s') = r + \gamma Q(s, \tau(s))$

- Given demonstration data $\langle s_1, r_1, d_1, s_2, r_2, d_2, \dots \rangle$
- Train Q(s,s') and au(s)

Cycle Loss



Model
$$au_{\psi}(s) = s'_{\tau}$$

Inverse dynamics $I_{\omega}(s, s'_{\tau}) = a$
Forward dynamics $f_{\phi}(s, a) = s'_{f}$

$$\mathcal{L}_{\psi} = -Q_{\theta}(s, C(s, \tau_{\psi}(s)) + \beta \| \tau_{\psi}(s) - C(s, \tau_{\psi}(s)) \|$$
$$y = r + \gamma \min_{i=1,2} Q_{\theta'_{i}}(s', C(s', \tau_{\psi'}(s')))$$



Model $au_{\psi}(s) = s'_{\tau}$ QSS $Q(s, s'_{\tau}) = q$ Forward dynamics $f_{\phi}(s, q) = s'_{f}$

$$\mathcal{L}_{\psi} = -Q_{\theta}(s, C(s, \tau_{\psi}(s)) + \beta \| \tau_{\psi}(s) - C(s, \tau_{\psi}(s)) \|$$

$$y = r + \gamma \min_{i=1,2} Q_{\theta'_{i}}(s', C(s', \tau_{\psi'}(s')))$$





 $s \to C(s, \tau(s)) \to \dots$



- 1) Select state with largest value
- 2) Determine action

$$s'_{\tau} \leftarrow \tau(s)$$



- 1) Select state with largest value
- 2) Determine action

$$s'_{\tau} \leftarrow \tau(s)$$



- 1) Select state with largest value
- 2) Determine action

$$\pi(s) \leftarrow I(s, s_{\tau}')$$

Table: Learning from observation with noise injected into policies

		Reacher-v2		InvertedPendulum-v2			
% Random	π_o	BCO	D3G	π_o	BCO	D3G	
0	-4.1 ± 0.7	-4.2 ± 0.6	-14.7 ± 30.5	1000 ± 0	1000 ± 0	3.0 ± 0.9	
25	-12.5 ± 1.0	-4.3 ± 0.6	-4.2 \pm 0.6	52.3 ± 3.7	1000 ± 0	602.1 ± 487.4	
50	-22.6 ± 0.9	-4.9 ± 0.7	-4.2 ± 0.6	18.0 ± 2.4	12.1 ± 8.3	900.2 ± 299.2	
75	-32.6 ± 0.4	$\textbf{-}6.6 \pm 1.3$	-4.6 ± 0.6	11.4 ± 1.3	12.1 ± 8.3	1000 ± 0	
100	-40.6 ± 0.5	-9.7 ± 0.8	-6.4 \pm 0.7	8.6 ± 0.3	31.0 ± 4.7	1000 ± 0	

Conclusion

- Introduced new value function, Q(s,s')
- Described predictive model for maximizing these values
- Showed benefits in redundant action spaces, transfer across permuted actions, and learning from observation
- Code: <u>https://github.com/uber-research/D3G</u>