

Deep Reinforcement Learning

Shan-Hung Wu
shwu@cs.nthu.edu.tw

Department of Computer Science,
National Tsing Hua University, Taiwan

Machine Learning

Outline

1 Introduction

2 Value-based Deep RL

- Deep Q -Network
- Improvements

3 Policy-based Deep RL

- Pathwise Derivative Methods
- Policy Gradient/Optimization Methods
- Variance Reduction and Actor-Critic

Outline

① Introduction

② Value-based Deep RL

- Deep Q -Network
- Improvements

③ Policy-based Deep RL

- Pathwise Derivative Methods
- Policy Gradient/Optimization Methods
- Variance Reduction and Actor-Critic

(Tabular) RL

- Q -learning:

$$Q^*(s, a) \leftarrow Q^*(s, a) + \eta [(R(s, a, s') + \gamma \max_{a'} Q^*(s', a')) - Q^*(s, a)]$$

- SARSA:

$$Q_\pi(s, a) \leftarrow Q_\pi(s, a) + \eta [(R(s, a, s') + \gamma Q_\pi(s', \pi(s')))] - Q_\pi(s, a)]$$

(Tabular) RL

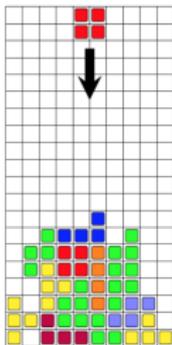
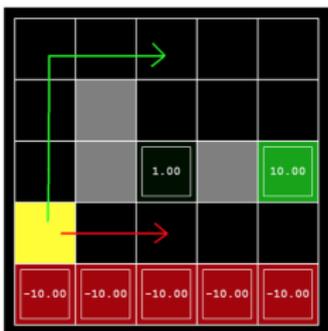
- Q -learning:

$$Q^*(s, a) \leftarrow Q^*(s, a) + \eta [(R(s, a, s') + \gamma \max_{a'} Q^*(s', a')) - Q^*(s, a)]$$

- SARSA:

$$Q_\pi(s, a) \leftarrow Q_\pi(s, a) + \eta [(R(s, a, s') + \gamma Q_\pi(s', \pi(s'))) - Q_\pi(s, a)]$$

- In realistic environments with large state/action space, requires a large table to store Q^*/Q_π values
 - Maze: $O(10^1)$, Tetris: $O(10^{60})$, Atari: $O(10^{16922})$ pixels
 - Continuous states/actions?
- May not be able to visit all (s, a) 's in limited training time



Generalizing across States

- Idea: to learn a function $f_{Q^*}(s, \mathbf{a}; \Theta)$ (resp. f_{Q_π}) that approximates $Q^*(s, \mathbf{a})$ (resp. $Q_\pi(s, \mathbf{a})$), $\forall s, \mathbf{a}$
 - Trained by a small number (millions) of samples
 - Generalizes to unseen states/actions
 - Smaller Θ to store

Generalizing across States

- Idea: to learn a function $f_{Q^*}(s, \mathbf{a}; \Theta)$ (resp. f_{Q_π}) that approximates $Q^*(s, \mathbf{a})$ (resp. $Q_\pi(s, \mathbf{a})$), $\forall s, \mathbf{a}$
 - Trained by a small number (millions) of samples
 - Generalizes to unseen states/actions
 - Smaller Θ to store
- E.g., in Q -learning, Q^* should satisfy Bellman optimality equation:

$$Q^*(s, \mathbf{a}) \leftarrow \sum_{s'} P(s'|s; \mathbf{a}) [R(s, \mathbf{a}, s') + \gamma \max_{a'} Q^*(s', \mathbf{a}')], \forall s, \mathbf{a}$$

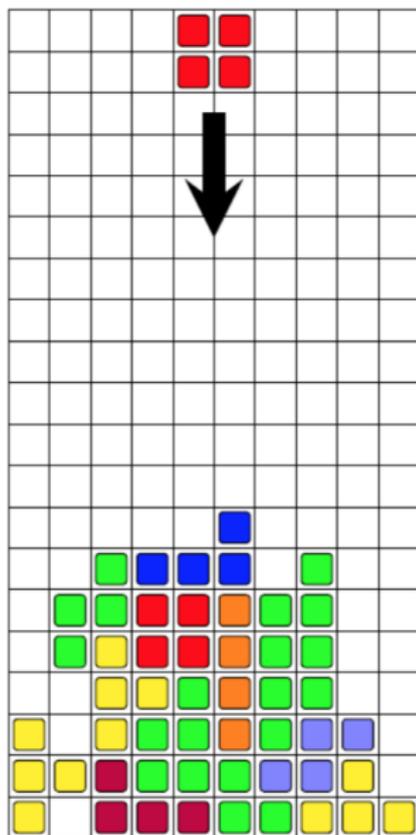
- Algorithm (TD estimate): initialize Θ arbitrarily, iterate until converge:
 - ① Take action \mathbf{a} from s using some exploration policy π' derived from f_{Q^*} (e.g., ϵ -greedy)
 - ② Observe s' and reward $R(s, \mathbf{a}, s')$, update Θ using SGD:

$$\Theta \leftarrow \Theta - \eta \nabla_{\Theta} C, \text{ where}$$

$$C(\Theta) = \left[R(s, \mathbf{a}, s') + \gamma \max_{a'} f_{Q^*}(s', \mathbf{a}'; \Theta) - f_{Q^*}(s, \mathbf{a}; \Theta) \right]^2$$

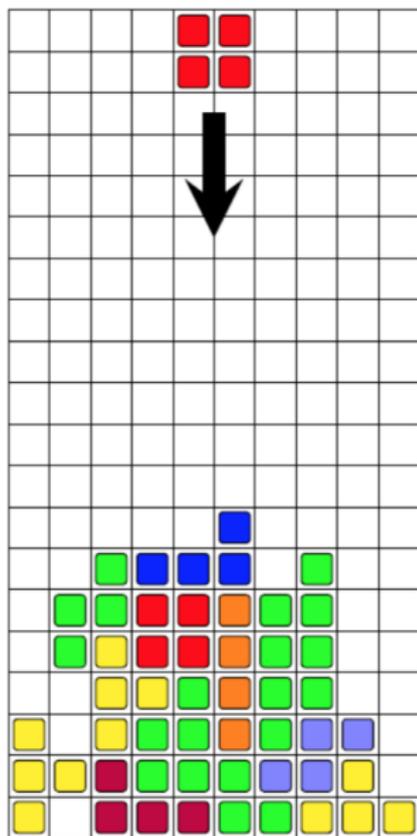
Works If with Careful Feature Engineering

- Tetris: [1]
 - States: $O(10^{60})$ configurations
 - Actions: rotation and translation to falling piece
- $f(s, \mathbf{a}; \Theta)$ and $C(\Theta)$ modeled as an approximated linear programming problem
- Hand-crafted features (22 in total)



Works If with Careful Feature Engineering

- Tetris: [1]
 - States: $O(10^{60})$ configurations
 - Actions: rotation and translation to falling piece
- $f(s, \mathbf{a}; \Theta)$ and $C(\Theta)$ modeled as an approximated linear programming problem
- Hand-crafted features (22 in total)
- Why not use a deep neural network to represent Q_Θ ?
 - One model for different tasks
 - Automatically learned features



Deep RL

- **Value-based:** use DNNs to represent **value/Q-function**
 - E.g., DQN
 - $\pi^*(s) \leftarrow \arg \max_a Q^*(s, a)$ only feasible if actions are discrete
- **Policy-based:** use DNNs to represent **policy π**
 - E.g., DDPG, Action-Critic, A3C, TRPO, PPO
- **Model-based:** deep RL when MDP/env. model is known
 - E.g., AlphaGo



Kohl and Stone, 2004



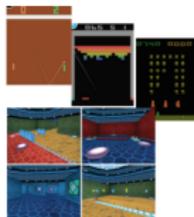
Ng et al, 2004



Tedrake et al, 2005



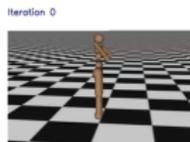
Kober and Peters, 2009



Mnih et al 2013 (DQN)
Mnih et al, 2015 (A3C)



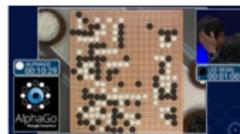
Silver et al, 2014 (DPG)
Lillicrap et al, 2015 (DDPG)



Schulman et al,
2016 (TRPO + GAE)



Levine*, Finn*, et
al, 2016
(GPS)



Silver*, Huang*, et
al, 2016
(AlphaGo)

Outline

① Introduction

② Value-based Deep RL

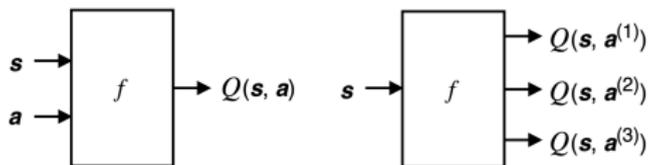
- Deep Q -Network
- Improvements

③ Policy-based Deep RL

- Pathwise Derivative Methods
- Policy Gradient/Optimization Methods
- Variance Reduction and Actor-Critic

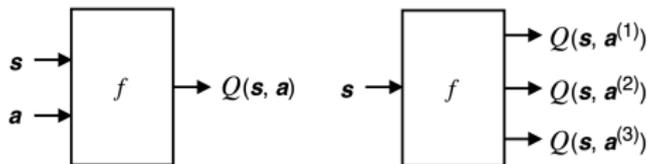
DNNs for Q^*

- Use a DNN $f_{Q^*}(s, a; \Theta)$ to represent $Q^*(s, a)$



DNNs for Q^*

- Use a DNN $f_{Q^*}(s, \mathbf{a}; \Theta)$ to represent $Q^*(s, \mathbf{a})$



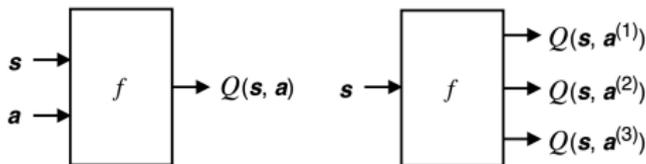
- Algorithm (TD): initialize Θ arbitrarily, iterate until converge:
 - Take action \mathbf{a} from s using some exploration policy π' derived from f_{Q^*} (e.g., ϵ -greedy)
 - Observe s' and reward $R(s, \mathbf{a}, s')$, update Θ using SGD:

$$\Theta \leftarrow \Theta - \eta \nabla_{\Theta} C, \text{ where}$$

$$C(\Theta) = \left[R(s, \mathbf{a}, s') + \gamma \max_{\mathbf{a}'} f_{Q^*}(s', \mathbf{a}'; \Theta) - f_{Q^*}(s, \mathbf{a}; \Theta) \right]^2$$

DNNs for Q^*

- Use a DNN $f_{Q^*}(s, \mathbf{a}; \Theta)$ to represent $Q^*(s, \mathbf{a})$



- Algorithm (TD): initialize Θ arbitrarily, iterate until converge:
 - Take action \mathbf{a} from s using some exploration policy π' derived from f_{Q^*} (e.g., ϵ -greedy)
 - Observe s' and reward $R(s, \mathbf{a}, s')$, update Θ using SGD:

$$\Theta \leftarrow \Theta - \eta \nabla_{\Theta} C, \text{ where}$$

$$C(\Theta) = \left[R(s, \mathbf{a}, s') + \gamma \max_{\mathbf{a}'} f_{Q^*}(s', \mathbf{a}'; \Theta) - f_{Q^*}(s, \mathbf{a}; \Theta) \right]^2$$

- However, **diverges** due to
 - Samples are correlated (violates i.i.d. assumption of training examples)
 - Non-stationary target ($f_{Q^*}(s', \mathbf{a}')$) changes as Θ is updated for current \mathbf{a}

Outline

① Introduction

② Value-based Deep RL

- Deep Q -Network
- Improvements

③ Policy-based Deep RL

- Pathwise Derivative Methods
- Policy Gradient/Optimization Methods
- Variance Reduction and Actor-Critic

Deep Q-Network (DQN)

- Naive TD algorithm diverges due to:
 - Samples are correlated
 - Non-stationary target

Deep Q-Network (DQN)

- Naive TD algorithm diverges due to:
 - Samples are correlated
 - Non-stationary target
- Stabilization techniques proposed by (Nature) DQN [5]:
 - *Experience replay*
 - *Delayed target network*

Experience Replay

- Use a replay memory \mathbb{D} to store recently seen transitions (s, a, r, s') 's
- Sample a mini-batch from \mathbb{D} and update Θ

Experience Replay

- Use a replay memory \mathbb{D} to store recently seen transitions (s, \mathbf{a}, r, s') 's
- Sample a mini-batch from \mathbb{D} and update Θ
- Algorithm (TD): initialize Θ arbitrarily, iterate until converge:
 - ① Take action \mathbf{a} from s using π' derived from f_{Q^*} (e.g., ϵ -greedy)
 - ② Observe s' and reward R , add (s, \mathbf{a}, R, s') to \mathbb{D}
 - ③ Sample a mini-batch of $(s^{(i)}, \mathbf{a}^{(i)}, R^{(i)}, s^{(i+1)})$'s from \mathbb{D} , do:

$$\Theta \leftarrow \Theta - \eta \nabla_{\Theta} C, \text{ where}$$

$$C(\Theta) = \sum_i \left[R^{(i)} + \gamma \max_{\mathbf{a}'} f_{Q^*}(s^{(i+1)}, \mathbf{a}'; \Theta) - f_{Q^*}(s^{(i)}, \mathbf{a}^{(i)}; \Theta) \right]^2$$

Delayed Target Network

- To avoid chasing a moving target, set the **target value** at network output parametrized by *old* Θ^-

Delayed Target Network

- To avoid chasing a moving target, set the **target value** at network output parametrized by **old** Θ^-
- Algorithm (TD): initialize Θ arbitrarily and $\Theta^- = \Theta$, iterate until converge:
 - ① Take action \mathbf{a} from s using π' derived from f_{Q^*} (e.g., ϵ -greedy)
 - ② Observe s' and reward R , add (s, \mathbf{a}, R, s') to \mathbb{D}
 - ③ Sample a mini-batch of $(s^{(i)}, \mathbf{a}^{(i)}, R^{(i)}, s^{(i+1)})$'s from \mathbb{D} , do:

$$\Theta \leftarrow \Theta - \eta \nabla_{\Theta} C, \text{ where}$$

$$C(\Theta) = \sum_i \left[R^{(i)} + \gamma \max_{\mathbf{a}'} f_{Q^*}(s^{(i+1)}, \mathbf{a}'; \Theta^-) - f_{Q^*}(s^{(i)}, \mathbf{a}^{(i)}; \Theta) \right]^2$$

- ④ Update $\Theta^- \leftarrow \Theta$ every K iterations

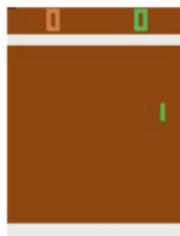
Other Tricks

- Optimization techniques matter in deep RL
 - Optimization error may lead to wrong traditions (trajectory)
 - And bad final policy

Other Tricks

- Optimization techniques matter in deep RL
 - Optimization error may lead to wrong traditions (trajectory)
 - And bad final policy
- *Reward clipping* for better conditioned gradients
 - Can't differentiate between small and large rewards
 - Better use batch normalization
- Use *RMSProp* instead of vanilla SGD for adaptive learning rate

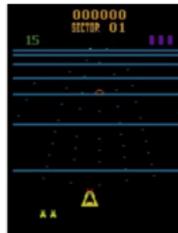
DQN on Atari



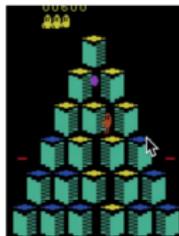
Pong



Enduro



Beamrider

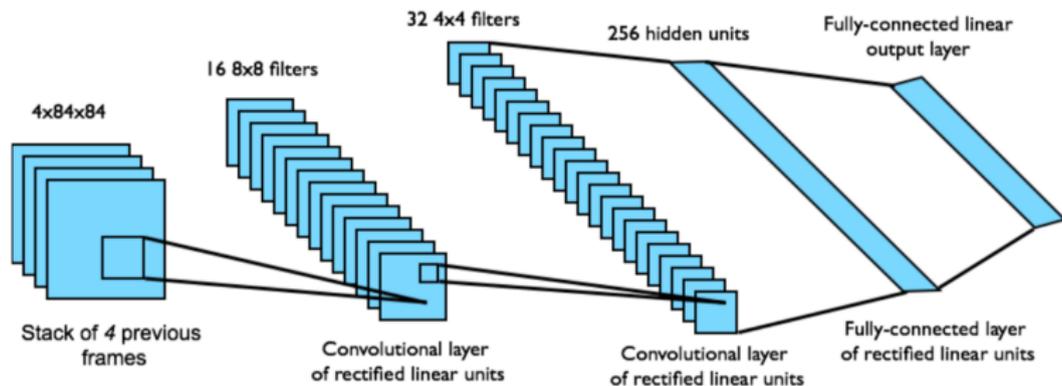


Q*bert

- 49 Atari 2600 games
- States: raw pixels
- Actions: 18 joystick/button positions
- Rewards: changes in score

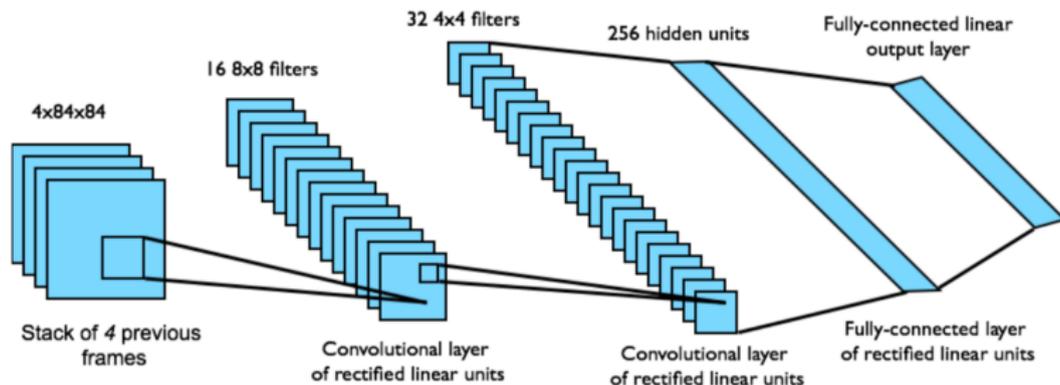
Network Architecture

- End-to-end from raw pixels to $Q^*(s, \mathbf{a})$
- CNN + fully connected layers
- Input: state s a stack of raw pixels from last 4 frames
- Output: 18 $Q^*(s, \mathbf{a})$'s (one for each action)



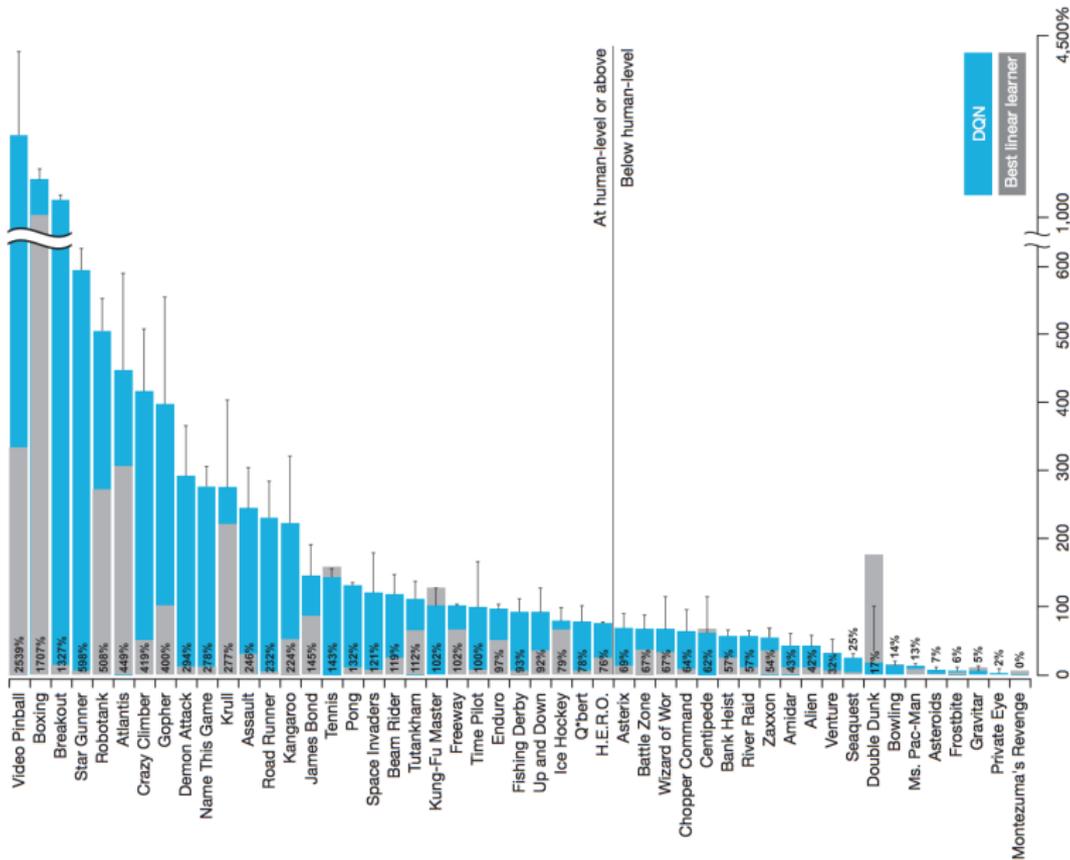
Network Architecture

- End-to-end from raw pixels to $Q^*(s, a)$
- CNN + fully connected layers
- Input: state s a stack of raw pixels from last 4 frames
- Output: 18 $Q^*(s, a)$'s (one for each action)



- Network architecture is *fixed across all games*

Results



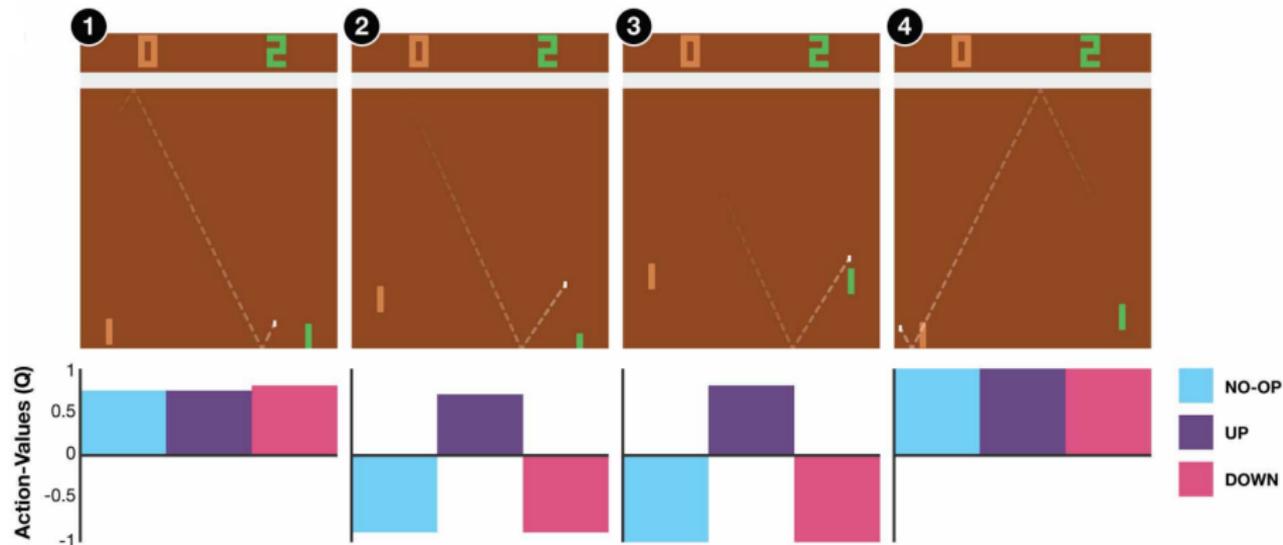
Effect of Stability Techniques

DQN

	Q-learning	Q-learning + Target Q	Q-learning + Replay	Q-learning + Replay + Target Q
Breakout	3	10	241	317
Enduro	29	142	831	1006
River Raid	1453	2868	4103	7447
Seaquest	276	1003	823	2894
Space Invaders	302	373	826	1089

- Delayed target network is less useful for large networks

Predicted Q^* Values for Pong



Outline

① Introduction

② Value-based Deep RL

- Deep Q -Network
- Improvements

③ Policy-based Deep RL

- Pathwise Derivative Methods
- Policy Gradient/Optimization Methods
- Variance Reduction and Actor-Critic

Improvements since DQN

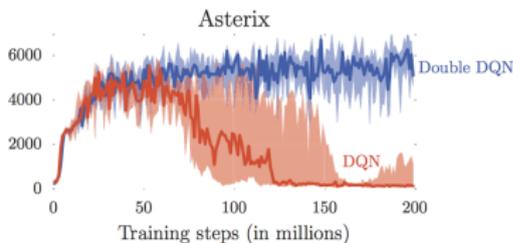
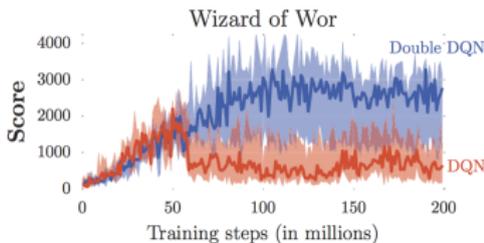
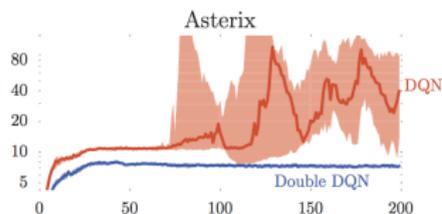
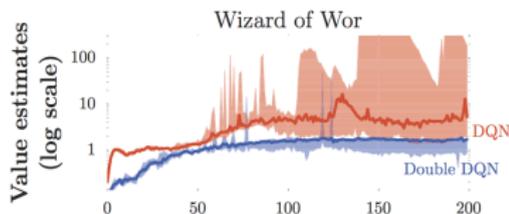
- Stabilization:
 - Double DQN [9]
 - Prioritized replay [7]
- Modeling additional prior:
 - Duelling network [10]
- Exploration:
 - NoisyNet [2]
- Large-scale implementation

Double DQN I

- DQN update rule: $\Theta \leftarrow \Theta - \eta \nabla_{\Theta} C$, where

$$C(\Theta) = \sum_i \left[R^{(i)} + \gamma \max_{a'} f_{Q^*}(s^{(i+1)}, a'; \Theta^-) - f_{Q^*}(s^{(i)}, a^{(i)}; \Theta) \right]^2$$

- There is an upward bias in $\max_{a'} f_{Q^*}(s^{(i+1)}, a'; \Theta^-)$
 - $f_{Q^*}(s^{(i+1)}, a'; \Theta^-)$ with high positive error is preferred
- At each step, the positive error is added to $f_{Q^*}(s^{(i)}, a^{(i)}; \Theta)$

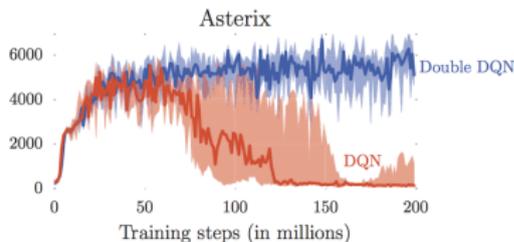
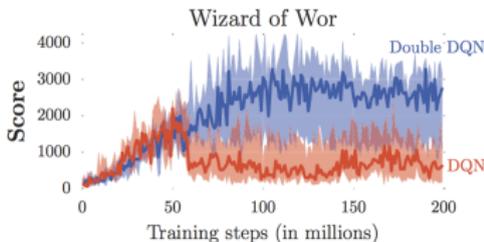
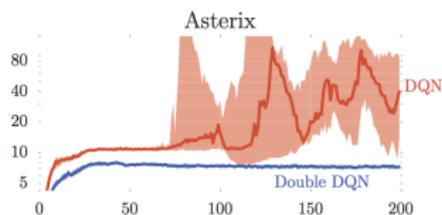
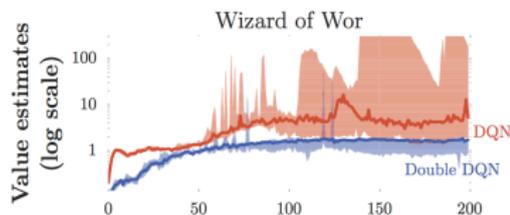


Double DQN II

- Double DQN (DDQN) [9]:

$$C(\Theta) = \sum_i \left[R^{(i)} + \gamma f_{Q^*}(s^{(i+1)}, \arg \max_{a'} f_{Q^*}(s^{(i+1)}, a'; \Theta); \Theta^-) - f_{Q^*}(s^{(i)}, a^{(i)}; \Theta) \right]^2$$

- Uses Θ to select the best action
- Uses Θ^- to evaluate the best action
- Random (unbiased) error added to $f_{Q^*}(s^{(i)}, a^{(i)}; \Theta)$ at each step

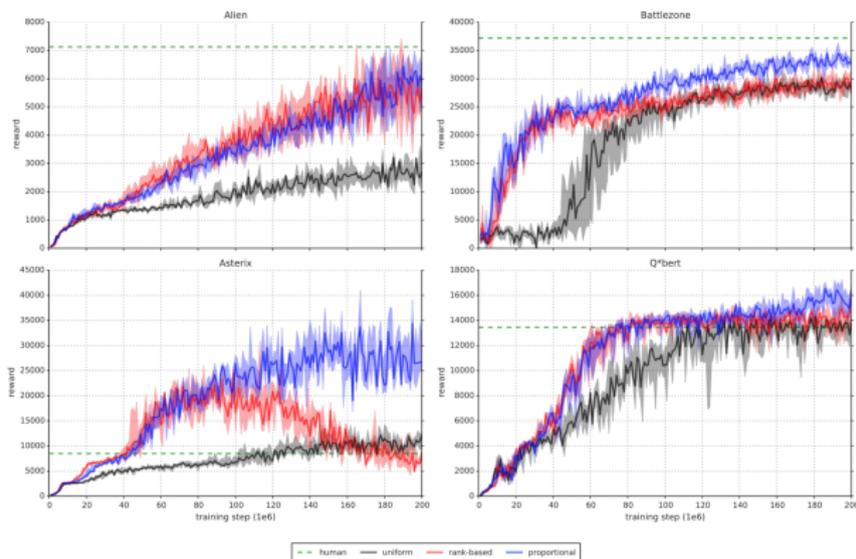


Prioritized Replay [7]

- Not all (s, \mathbf{a}, R, s') 's from \mathbb{D} are equally helpful to training f_{Q^*}
- Sample (s, \mathbf{a}, R, s') 's with probability proportional to “surprise” in terms of Bellman equation:

$$|R + \gamma \max_{\mathbf{a}'} f_{Q^*}(s', \mathbf{a}'; \Theta^-) - f_{Q^*}(s, \mathbf{a}; \Theta)|$$

- Rank-based alternative: \mathbb{D} a priority queue

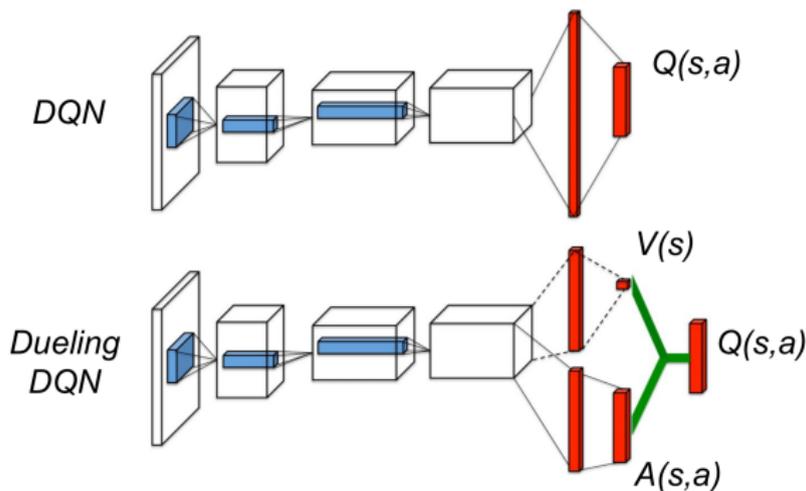


Dueling Network [10]

- $Q^*(s, a) = V^*(s) + A^*(s, a)$
 - $A^*(s, a)$ the *advantage function* of a

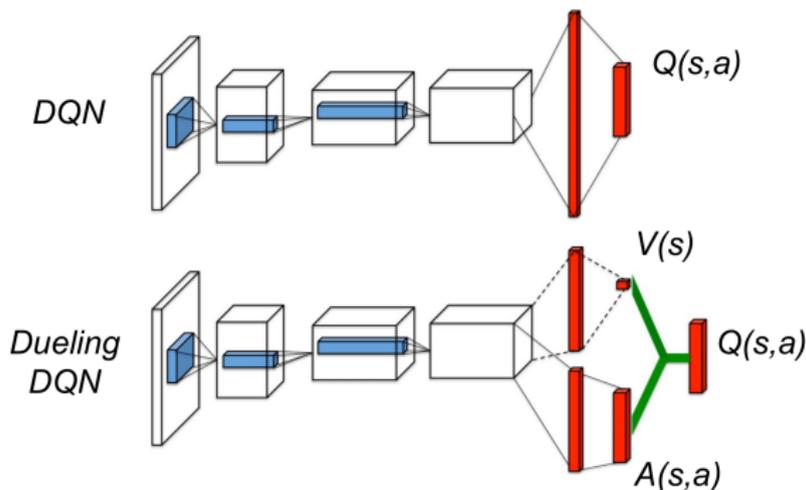
Dueling Network [10]

- $Q^*(s, a) = V^*(s) + A^*(s, a)$
 - $A^*(s, a)$ the **advantage function** of a
- Idea: to model this prior and learn $f_{Q^*}(s, a) = f_{V^*}(s) + f_{A^*}(s, a)$



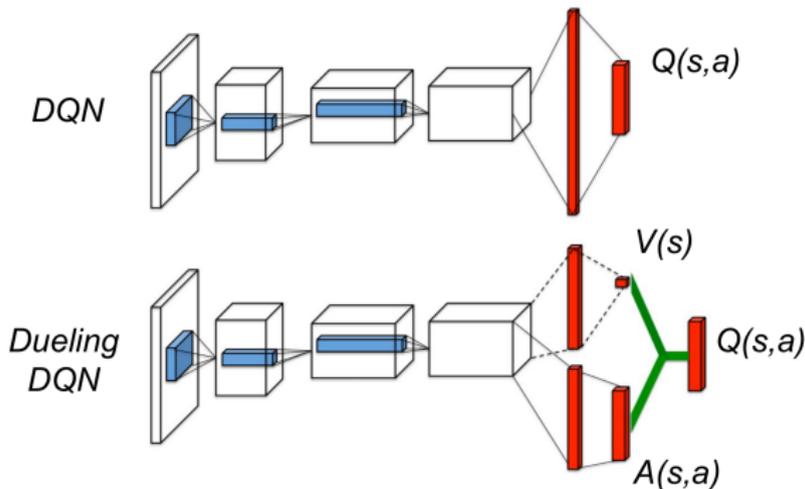
Dueling Network [10]

- $Q^*(s, a) = V^*(s) + A^*(s, a)$
 - $A^*(s, a)$ the **advantage function** of a
- Idea: to model this prior and learn $f_{Q^*}(s, a) = f_{V^*}(s) + f_{A^*}(s, a)$
 - Not well-defined: $f_{Q^*}(s, a) = (f_{V^*}(s) + c) + (f_{A^*}(s, a) - c)$ for any c
- Dueling DQN: $f_{Q^*}(s, a) = f_{V^*}(s) + (f_{A^*}(s, a) - \max_{a'} f_{A^*}(s, a'))$
 - The best action a^* has zero advantage and $f_{Q^*}(s, a^*) = f_{V^*}(s)$

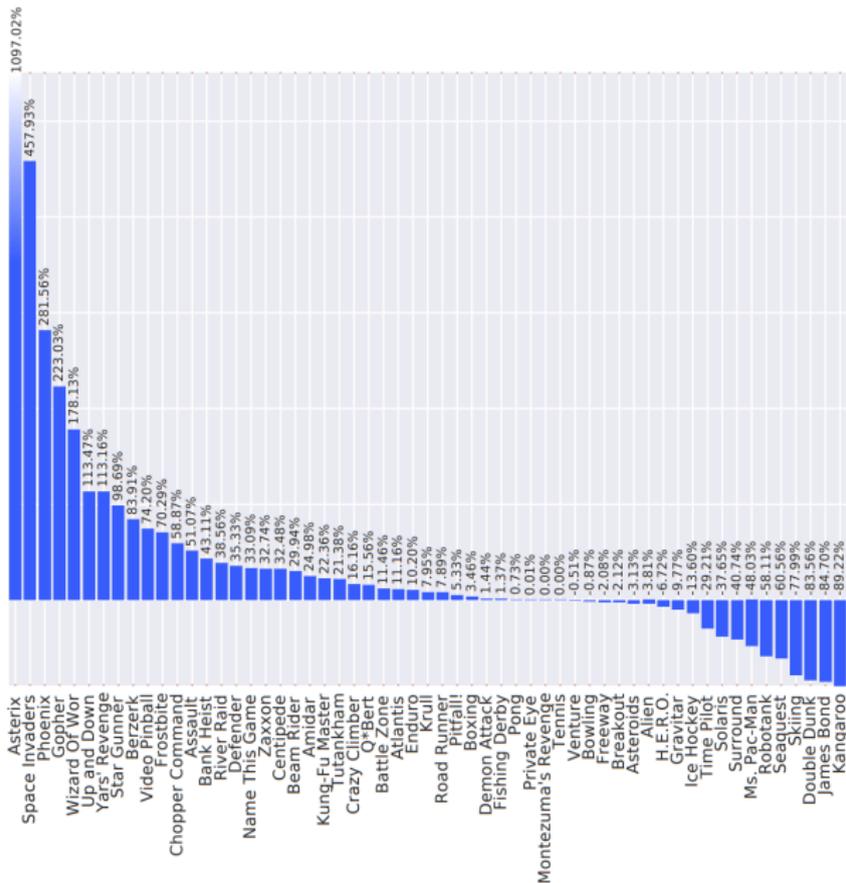


Dueling Network [10]

- $Q^*(s, a) = V^*(s) + A^*(s, a)$
 - $A^*(s, a)$ the **advantage function** of a
- Idea: to model this prior and learn $f_{Q^*}(s, a) = f_{V^*}(s) + f_{A^*}(s, a)$
 - Not well-defined: $f_{Q^*}(s, a) = (f_{V^*}(s) + c) + (f_{A^*}(s, a) - c)$ for any c
- Dueling DQN: $f_{Q^*}(s, a) = f_{V^*}(s) + (f_{A^*}(s, a) - \max_{a'} f_{A^*}(s, a'))$
 - The best action a^* has zero advantage and $f_{Q^*}(s, a^*) = f_{V^*}(s)$
- Stabilized version: $f_{Q^*}(s, a) = f_{V^*}(s) + (f_{A^*}(s, a) - \frac{1}{|\mathbb{A}|} \sum_{a'} f_{A^*}(s, a'))$
 - f_{V^*} and f_{A^*} are off-target (by a constant) but f_{A^*} changes more slowly

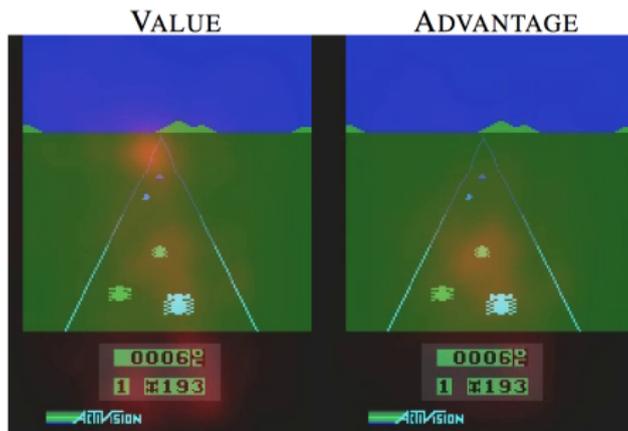


Improvement over Prioritized DDQN



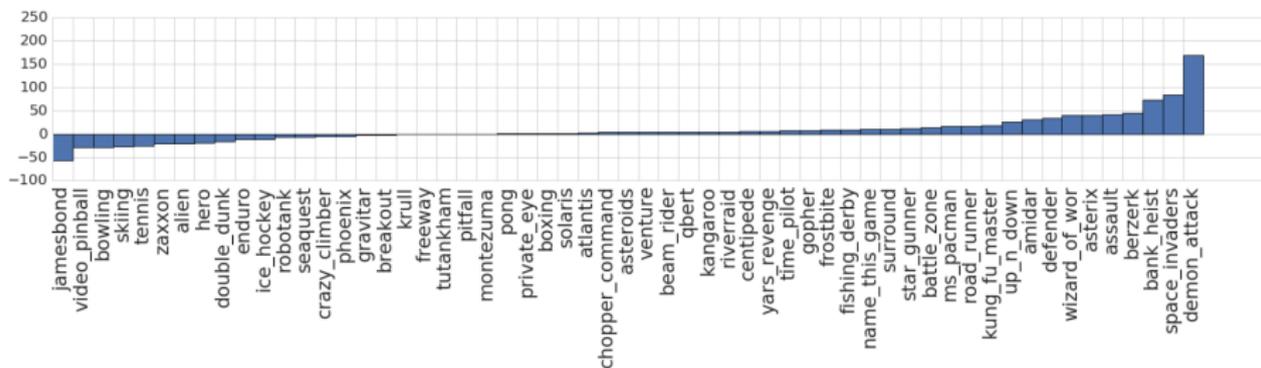
See, Attend, and Drive

- Atari game: Enduro
- Attention mask: $\frac{\partial f_{V^*}}{\partial s}(s)$ and $\frac{\partial f_{A^*}}{\partial s}(s)$
- $f_{Q^*}(s, \mathbf{a}) = f_{V^*}(s) + f_{A^*}(s, \mathbf{a})$
 - f_{V^*} pays attention to the road
 - f_{A^*} pays attention only when there's obstacles in front



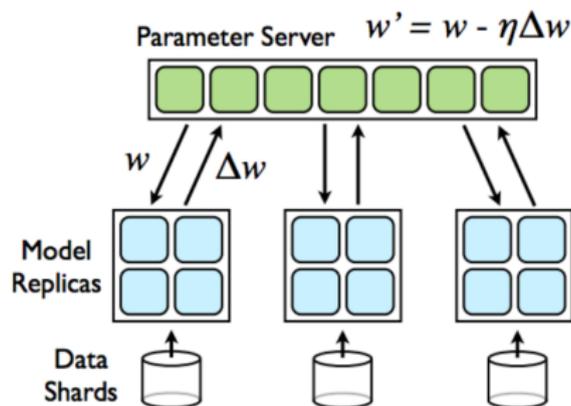
NoisyNet [2]

- Instead of using ϵ -greedy for exploration, add noise to Θ
- The level of noise is learned by SGD along with Θ
- Improvement over Dueling DQN:



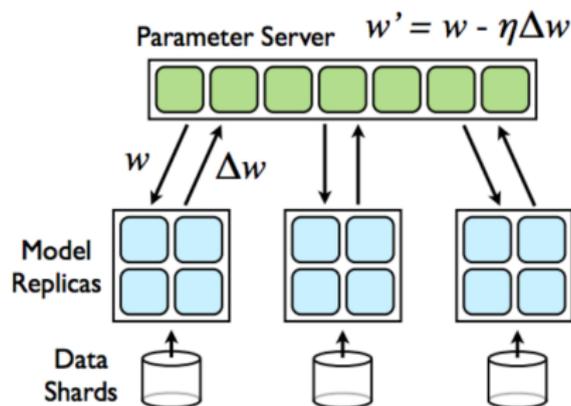
Scaling Up DQN on Single Machine

- Exploits multi-threading of modern CPUs/GPUs
- Run/train multiple agents in parallel (one per thread/GPU)
 - Θ shared between threads in main memory
 - Data-parallelism



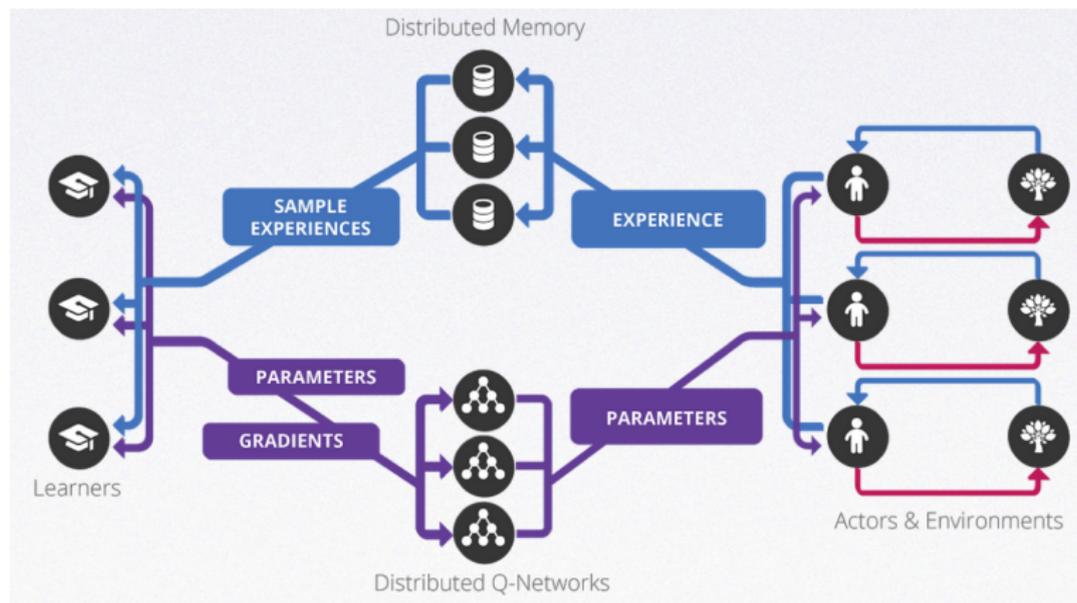
Scaling Up DQN on Single Machine

- Exploits multi-threading of modern CPUs/GPUs
- Run/train multiple agents in parallel (one per thread/GPU)
 - Θ shared between threads in main memory
 - Data-parallelism
- Parallelism *decorrelates samples*
 - Alternative to experience replay



Scaling Out DQN with Gorila [6]

- Distributed system architecture for large-scale RL
- 10x faster than Nature DQN
- Applied to recommender systems in Google



Outline

① Introduction

② Value-based Deep RL

- Deep Q -Network
- Improvements

③ Policy-based Deep RL

- Pathwise Derivative Methods
- Policy Gradient/Optimization Methods
- Variance Reduction and Actor-Critic

Why Policy Network?

- Policy-based deep RL: use a DNN $g_{\pi}(s; \Phi)$ to approximate $\pi(s)$
- Why?

Why Policy Network?

- Policy-based deep RL: use a DNN $g_{\pi}(s; \Phi)$ to approximate $\pi(s)$
- Why?
- In DQN (or any method based on value/policy iteration), one needs to solve

$$\pi^* = \arg \max_{a'} Q^*(s, a') \text{ or } \hat{\pi} = \arg \max_{a'} Q^{\pi}(s, a')$$

- Not applicable to continuous action space \mathbb{A} common in, e.g., robotics



Why Policy Network?

- Policy-based deep RL: use a DNN $g_{\pi}(s; \Phi)$ to approximate $\pi(s)$
- Why?
- In DQN (or any method based on value/policy iteration), one needs to solve

$$\pi^* = \arg \max_{a'} Q^*(s, a') \text{ or } \hat{\pi} = \arg \max_{a'} Q^{\pi}(s, a')$$

- Not applicable to continuous action space \mathbb{A} common in, e.g., robotics
- π may be easier to learn than Q or V



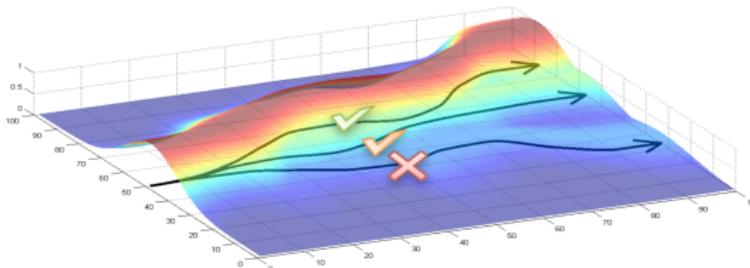
Modeling π

- π can be either
 - deterministic: $g_{\pi}(\mathbf{s}; \Phi) = \mathbf{a}$, or
 - stochastic: $g_{\pi}(\mathbf{s}; \Phi) = \mathbf{P}(\mathbf{a}|\mathbf{s})$
- ***Pathwise derivative methods***
 - For deterministic π and continuous \mathbb{A}

- ***Policy gradient/optimization*** methods
 - For stochastic π

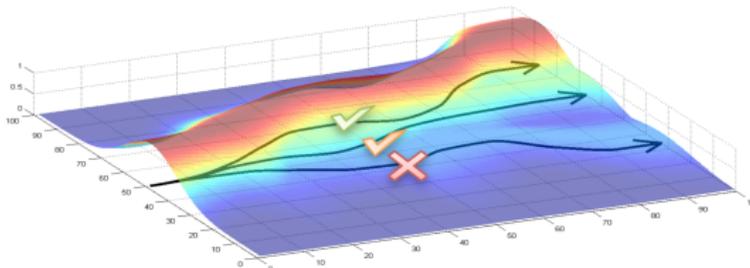
Modeling π

- π can be either
 - deterministic: $g_\pi(s; \Phi) = \mathbf{a}$, or
 - stochastic: $g_\pi(s; \Phi) = \mathbf{P}(\mathbf{a}|s)$
- **Pathwise derivative methods**
 - For deterministic π and continuous \mathbb{A}
 - To find Φ such that $g_\pi(s; \Phi)$ gives action \mathbf{a} maximizing $Q^*(s, \mathbf{a})$
 - Changes the trajectory of an episode in the graph of accumulative rewards
- **Policy gradient/optimization** methods
 - For stochastic π



Modeling π

- π can be either
 - deterministic: $g_{\pi}(s; \Phi) = \mathbf{a}$, or
 - stochastic: $g_{\pi}(s; \Phi) = P(\mathbf{a}|s)$
- **Pathwise derivative methods**
 - For deterministic π and continuous \mathbb{A}
 - To find Φ such that $g_{\pi}(s; \Phi)$ gives action \mathbf{a} maximizing $Q^*(s, \mathbf{a})$
 - Changes the trajectory of an episode in the graph of accumulative rewards
- **Policy gradient/optimization** methods
 - For stochastic π
 - To find Φ such that gives trajectory of high accumulative rewards
 - Do **not** change the trajectory (but its probability) of an episode



Outline

① Introduction

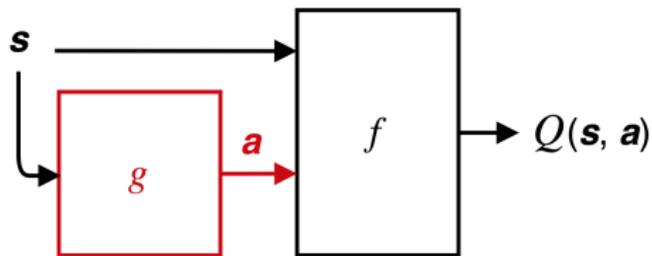
② Value-based Deep RL

- Deep Q -Network
- Improvements

③ Policy-based Deep RL

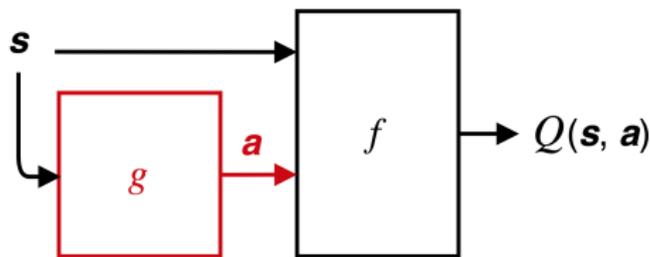
- Pathwise Derivative Methods
- Policy Gradient/Optimization Methods
- Variance Reduction and Actor-Critic

Deep Deterministic Policy Gradient (DDPG) [4]



- Based on DQN
 - Q -learning is off-policy and works with changing exploration strategies
- Deterministic policy: $g_{\pi^*}(s; \Phi) = a \in \mathbb{R}$

Deep Deterministic Policy Gradient (DDPG) [4]



- Based on DQN
 - Q -learning is off-policy and works with changing exploration strategies
- Deterministic policy: $g_{\pi^*}(s; \Phi) = \mathbf{a} \in \mathbb{R}$
- Goal: to find Φ maximizing $E_s [f_{Q^*}(s, \mathbf{a}; \Theta)]$, where $\mathbf{a} = g_{\pi^*}(s; \Phi)$
- SGD update rule:

$$\begin{aligned}\Phi &\leftarrow \Phi + \eta \frac{\partial E_s [f_{Q^*}(s, \mathbf{a}; \Theta)]}{\partial \Phi} \\ &= \Phi + \eta E_s \left[\frac{\partial f_{Q^*}}{\partial \mathbf{a}}(s, \mathbf{a}; \Theta) \cdot \frac{\partial g_{\pi^*}}{\partial \Phi}(s; \Phi) \right]\end{aligned}$$

DDPG Algorithm (TD)

- Initialize Θ and Φ arbitrarily, set $\Theta^- = \Theta$ and $\Phi^- = \Phi$, iterate until converge:

- Take action $\mathbf{a} = g_{\pi^*}(s; \Phi) + \mathbf{z}$ from s , where \mathbf{z} is a random noise for exploration
- Observe s' and reward R , add (s, \mathbf{a}, R, s') to \mathbb{D}
- Sample a mini-batch of $(s^{(i)}, \mathbf{a}^{(i)}, R^{(i)}, s^{(i+1)})$'s from \mathbb{D}
- Update Θ :

$$\Theta \leftarrow \Theta - \eta \nabla_{\Theta} C, \text{ where}$$

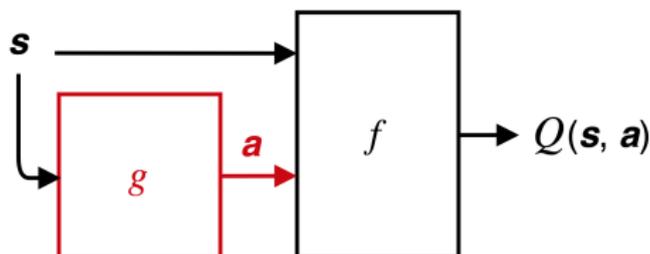
$$C(\Theta) = \sum_i \left[R^{(i)} + \gamma f_{Q^*}(s^{(i+1)}, g_{\pi^*}(s^{(i)}; \Phi^-); \Theta^-) - f_{Q^*}(s^{(i)}, \mathbf{a}^{(i)}; \Theta) \right]^2$$

- Update Φ :

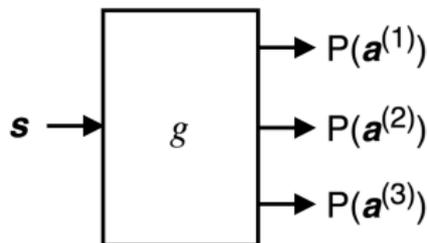
$$\Phi \leftarrow \Phi + \lambda \sum_i \frac{\partial f_{Q^*}}{\partial \mathbf{a}}(s^{(i)}, g_{\pi^*}(s^{(i)}; \Phi); \Theta) \cdot \frac{\partial g_{\pi^*}}{\partial \Phi}(s^{(i)}; \Phi)$$

- Update $\Theta^- \leftarrow \tau \Theta + (1 - \tau) \Theta^-$ and $\Phi^- \leftarrow \tau \Phi + (1 - \tau) \Phi^-$

Limitations



- Only applicable to continuous action space \mathbb{A}
- Cannot backprop through samples when calculating $\frac{\partial E_{\mathbf{s}}[f_{Q^*}(\mathbf{s}; g_{\pi^*}(\mathbf{s}; \Phi); \Theta)]}{\partial \Phi}$
- For discrete \mathbb{A} , it's more natural to use a DNN to model a stochastic policy: $g_{\pi}(\mathbf{s}) = P(\mathbf{a}|\mathbf{s}), \forall \mathbf{a}$



Outline

① Introduction

② Value-based Deep RL

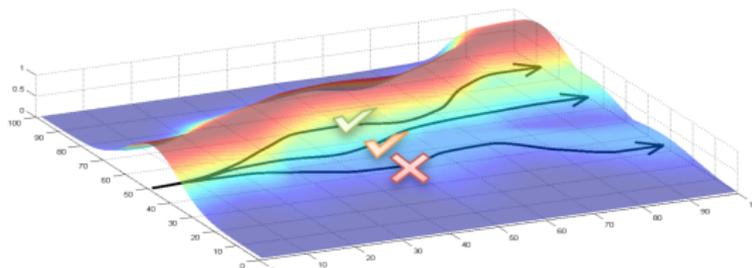
- Deep Q -Network
- Improvements

③ Policy-based Deep RL

- Pathwise Derivative Methods
- Policy Gradient/Optimization Methods
- Variance Reduction and Actor-Critic

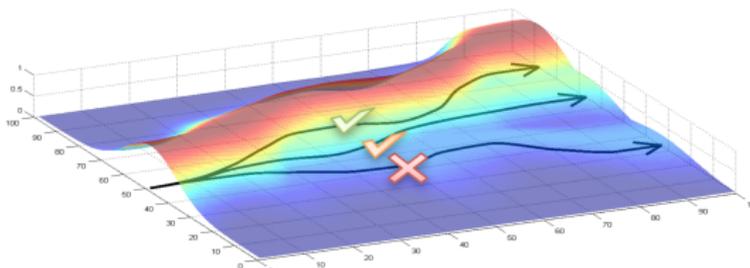
Episodic Policy Gradient

- **Policy gradient/optimization** methods
 - For stochastic policy: $g_{\pi}(s) = P(\mathbf{a}|s; \Phi), \forall \mathbf{a}$ (discrete or continuous)
 - Do **not** change the trajectory (but its probability) of an episode



Episodic Policy Gradient

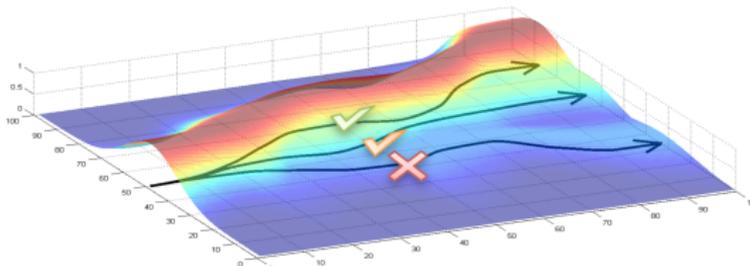
- **Policy gradient/optimization** methods
 - For stochastic policy: $g_{\pi}(s) = P(\mathbf{a}|s; \Phi), \forall \mathbf{a}$ (discrete or continuous)
 - Do **not** change the trajectory (but its probability) of an episode
- Given an episode, let $\tau = \{(s^{(t)}, \mathbf{a}^{(t)}, R^{(t)}, s^{(t+1)})\}_t$ be the sequence of state-action transitions
 - Action $\mathbf{a}^{(t)}$ sampled from $g_{\pi}(s^{(t)})$



Episodic Policy Gradient

- **Policy gradient/optimization** methods
 - For stochastic policy: $g_{\pi}(s) = P(\mathbf{a}|s; \Phi), \forall \mathbf{a}$ (discrete or continuous)
 - Do **not** change the trajectory (but its probability) of an episode
- Given an episode, let $\tau = \{(s^{(t)}, \mathbf{a}^{(t)}, R^{(t)}, s^{(t+1)})\}_t$ be the sequence of state-action transitions
 - Action $\mathbf{a}^{(t)}$ sampled from $g_{\pi}(s^{(t)})$
- Let $R(\tau) = \sum_t \gamma^t R^{(t)}$, our goal:

$$\arg \max_{\Phi} E_{\tau} [R(\tau); \Phi] = \arg \max_{\Phi} \sum_{\tau} P(\tau; \Phi) R(\tau)$$



Policy Gradient

- Let $J(\Phi) = \sum_{\tau} P(\tau; \Phi)R(\tau)$, we have:

$$\begin{aligned}\nabla_{\Phi} J(\Phi) &= \nabla_{\Phi} \sum_{\tau} P(\tau; \Phi)R(\tau) = \sum_{\tau} \nabla_{\Phi} P(\tau; \Phi)R(\tau) \\ &= \sum_{\tau} P(\tau; \Phi) \frac{\nabla_{\Phi} P(\tau; \Phi)}{P(\tau; \Phi)} R(\tau) \\ &= \sum_{\tau} P(\tau; \Phi) \nabla_{\Phi} \log P(\tau; \Phi) R(\tau)\end{aligned}$$

Policy Gradient

- Let $J(\Phi) = \sum_{\tau} P(\tau; \Phi) R(\tau)$, we have:

$$\begin{aligned}\nabla_{\Phi} J(\Phi) &= \nabla_{\Phi} \sum_{\tau} P(\tau; \Phi) R(\tau) = \sum_{\tau} \nabla_{\Phi} P(\tau; \Phi) R(\tau) \\ &= \sum_{\tau} P(\tau; \Phi) \frac{\nabla_{\Phi} P(\tau; \Phi)}{P(\tau; \Phi)} R(\tau) \\ &= \sum_{\tau} P(\tau; \Phi) \nabla_{\Phi} \log P(\tau; \Phi) R(\tau) \\ &= \sum_{\tau} P(\tau; \Phi) \nabla_{\Phi} \log \prod_t P(s^{(t+1)} | s^{(t)}, \mathbf{a}^{(t)}) P(\mathbf{a}^{(t)} | s^{(t)}; \Phi) R(\tau) \\ &= \sum_{\tau} P(\tau; \Phi) \nabla_{\Phi} \sum_t \left[\log P(s^{(t+1)} | s^{(t)}, \mathbf{a}^{(t)}) + \log P(\mathbf{a}^{(t)} | s^{(t)}; \Phi) \right] R(\tau) \\ &= \sum_{\tau} P(\tau; \Phi) \sum_t \nabla_{\Phi} \log P(\mathbf{a}^{(t)} | s^{(t)}; \Phi) R(\tau)\end{aligned}$$

Policy Gradient

- Let $J(\Phi) = \sum_{\tau} P(\tau; \Phi) R(\tau)$, we have:

$$\begin{aligned}\nabla_{\Phi} J(\Phi) &= \nabla_{\Phi} \sum_{\tau} P(\tau; \Phi) R(\tau) = \sum_{\tau} \nabla_{\Phi} P(\tau; \Phi) R(\tau) \\ &= \sum_{\tau} P(\tau; \Phi) \frac{\nabla_{\Phi} P(\tau; \Phi)}{P(\tau; \Phi)} R(\tau) \\ &= \sum_{\tau} P(\tau; \Phi) \nabla_{\Phi} \log P(\tau; \Phi) R(\tau) \\ &= \sum_{\tau} P(\tau; \Phi) \nabla_{\Phi} \log \prod_t P(s^{(t+1)} | s^{(t)}, \mathbf{a}^{(t)}) P(\mathbf{a}^{(t)} | s^{(t)}; \Phi) R(\tau) \\ &= \sum_{\tau} P(\tau; \Phi) \nabla_{\Phi} \sum_t \left[\log P(s^{(t+1)} | s^{(t)}, \mathbf{a}^{(t)}) + \log P(\mathbf{a}^{(t)} | s^{(t)}; \Phi) \right] R(\tau) \\ &= \sum_{\tau} P(\tau; \Phi) \sum_t \nabla_{\Phi} \log P(\mathbf{a}^{(t)} | s^{(t)}; \Phi) R(\tau) \\ &= \sum_{\tau} P(\tau; \Phi) \sum_t \nabla_{\Phi} \log P(\mathbf{a}^{(t)} | s^{(t)}; \Phi) \sum_{t'=t}^H \gamma^{t'} R^{(t')}(s^{(t)}, \mathbf{a}^{(t)}, s^{(t+1)}) \\ &= \sum_{\tau} P(\tau; \Phi) \sum_t \nabla_{\Phi} \log P(\mathbf{a}^{(t)} | s^{(t)}; \Phi) \sum_{t'=t}^H \gamma^{t'} R^{(t')}\end{aligned}$$

- Assumes that the environment is MDP-alike
 - But no need for the exact model

REINFORCE Algorithm

$$\nabla_{\Phi} J(\Phi) = \sum_{\tau} P(\tau; \Phi) \sum_t \nabla_{\Phi} \log P(\mathbf{a}^{(t)} | \mathbf{s}^{(t)}; \Phi) \sum_{t'=t}^H \gamma^{t'} R^{(t')}$$

- REINFORCE (MC estimate): initialize Φ arbitrarily, iterate until converge:
 - ① Run episodes $\{\tau^{(i)}\}_i$ by sampling actions from $g(\cdot; \Phi)$
 - ② For each time step t in an episode, compute $R^{(i,t)} = \sum_{t'=t}^{H^{(i)}} \gamma^{t'} R^{(i,t')}$
 - ③ Update Φ using SGD:

$$\Phi \leftarrow \Phi + \eta \nabla_{\Phi} \hat{J}, \text{ where}$$

$$\nabla_{\Phi} \hat{J}(\Phi) = \sum_{i,t} \nabla_{\Phi} \log P(\mathbf{a}^{(i,t)} | \mathbf{s}^{(i,t)}; \Phi) R^{(i,t)}.$$

REINFORCE Algorithm

$$\nabla_{\Phi} J(\Phi) = \sum_{\tau} P(\tau; \Phi) \sum_t \nabla_{\Phi} \log P(\mathbf{a}^{(t)} | \mathbf{s}^{(t)}; \Phi) \sum_{t'=t}^H \gamma^{t'} R^{(t')}$$

- REINFORCE (MC estimate): initialize Φ arbitrarily, iterate until converge:
 - ① Run episodes $\{\tau^{(i)}\}_i$ by sampling actions from $g(\cdot; \Phi)$
 - ② For each time step t in an episode, compute $R^{(i,t)} = \sum_{t'=t}^{H^{(i)}} \gamma^{t'} R^{(i,t')}$
 - ③ Update Φ using SGD:

$$\Phi \leftarrow \Phi + \eta \nabla_{\Phi} \hat{J}, \text{ where}$$

$$\nabla_{\Phi} \hat{J}(\Phi) = \sum_{i,t} \nabla_{\Phi} \log P(\mathbf{a}^{(i,t)} | \mathbf{s}^{(i,t)}; \Phi) R^{(i,t)}.$$

- REINFORCE-style policy gradient: $\nabla \log \text{prob. of actions} \times \text{episodic rewards}$

Outline

① Introduction

② Value-based Deep RL

- Deep Q -Network
- Improvements

③ Policy-based Deep RL

- Pathwise Derivative Methods
- Policy Gradient/Optimization Methods
- Variance Reduction and Actor-Critic

Variance

$$\nabla_{\Phi} J(\Phi) \propto \sum_{i,t} \nabla_{\Phi} \log P(\mathbf{a}^{(i,t)} | \mathbf{s}^{(i,t)}; \Phi) \sum_{t'=t}^{H^{(i)}} \gamma^{t'} R^{(i,t')}$$

- $\sum_{t'=t}^{H^{(i)}} \gamma^{t'} R^{(i,t')}$ is an MC estimate of $Q_{\pi}(\mathbf{s}^{(i,t)}, \mathbf{a}^{(i,t)}) = \mathbb{E}_{\{\mathbf{s}^{(t')}, \mathbf{a}^{(t')}\}_{t'}} \left[\sum_{t'} \gamma^{t'} R^{(t')} | \mathbf{s}^{(0)} = \mathbf{s}^{(i,t)}, \mathbf{a}^{(0)} = \mathbf{a}^{(i,t)} \right]$
 - using samples rolled out from single episode

Variance

$$\nabla_{\Phi} J(\Phi) \propto \sum_{i,t} \nabla_{\Phi} \log P(\mathbf{a}^{(i,t)} | \mathbf{s}^{(i,t)}; \Phi) \sum_{t'=t}^{H(i)} \gamma^{t'} R^{(i,t')}$$

- $\sum_{t'=t}^{H(i)} \gamma^{t'} R^{(i,t')}$ is an MC estimate of $Q_{\pi}(\mathbf{s}^{(i,t)}, \mathbf{a}^{(i,t)}) = \mathbb{E}_{\{\mathbf{s}^{(t')}, \mathbf{a}^{(t')}\}_{t'}} \left[\sum_{t'} \gamma^{t'} R^{(t')} | \mathbf{s}^{(0)} = \mathbf{s}^{(i,t)}, \mathbf{a}^{(0)} = \mathbf{a}^{(i,t)} \right]$
 - using samples rolled out from single episode
- TD vs. MC estimate:
 - TD: biased, but low variance
 - MC: unbiased, but **high variance**
- How to lower the variance of vanilla policy gradient algorithm?

Variance

$$\nabla_{\Phi} J(\Phi) \propto \sum_{i,t} \nabla_{\Phi} \log P(\mathbf{a}^{(i,t)} | \mathbf{s}^{(i,t)}; \Phi) \sum_{t'=t}^{H^{(i)}} \gamma^{t'} R^{(i,t')}$$

- $\sum_{t'=t}^{H^{(i)}} \gamma^{t'} R^{(i,t')}$ is an MC estimate of $Q_{\pi}(\mathbf{s}^{(i,t)}, \mathbf{a}^{(i,t)}) = \mathbb{E}_{\{\mathbf{s}^{(t')}, \mathbf{a}^{(t')}\}_{t'}} \left[\sum_{t'} \gamma^{t'} R^{(t')} | \mathbf{s}^{(0)} = \mathbf{s}^{(i,t)}, \mathbf{a}^{(0)} = \mathbf{a}^{(i,t)} \right]$
 - using samples rolled out from single episode
- TD vs. MC estimate:
 - TD: biased, but low variance
 - MC: unbiased, but **high variance**
- How to lower the variance of vanilla policy gradient algorithm?
 - To reduce the magnitude of $\sum_{t'=t}^{H^{(i)}} \gamma^{t'} R^{(i,t')}$
 - Eg., use a smaller γ or subtract a **baseline** from $\sum_{t'=t}^{H^{(i)}} \gamma^{t'} R^{(i,t')}$
 - To approximate $\sum_{t'=t}^{H^{(i)}} \gamma^{t'} R^{(i,t')}$ by a DNN and take advantage of its generalizability
 - To collect more samples

Baseline I

$$\nabla_{\Phi} J(\Phi) \propto \sum_{i,t} \nabla_{\Phi} \log P(\mathbf{a}^{(i,t)} | \mathbf{s}^{(i,t)}; \Phi) \left(\sum_{t'=t}^{H(i)} \gamma^{t'} R^{(i,t')} - b \right)$$

- b reduces variance without adding bias *as long as it's independent with actions*:

$$\begin{aligned} & \sum_{\tau} P(\tau; \Phi) \sum_t \nabla_{\Phi} \log P(\mathbf{a}^{(t)} | \mathbf{s}^{(t)}; \Phi) b \\ &= \sum_{\tau} P(\tau; \Phi) \nabla_{\Phi} \log P(\tau; \Phi) b \\ &= \sum_{\tau} P(\tau; \Phi) \frac{\nabla_{\Phi} P(\tau; \Phi)}{P(\tau; \Phi)} b \\ &= \nabla_{\Phi} \sum_{\tau} P(\tau; \Phi) b = \nabla_{\Phi} b = 0 \end{aligned}$$

Baseline I

$$\nabla_{\Phi} J(\Phi) \propto \sum_{i,t} \nabla_{\Phi} \log P(\mathbf{a}^{(i,t)} | \mathbf{s}^{(i,t)}; \Phi) \left(\sum_{t'=t}^{H(i)} \gamma^{t'} R^{(i,t')} - b \right)$$

- b reduces variance without adding bias *as long as it's independent with actions*:

$$\begin{aligned} & \sum_{\tau} P(\tau; \Phi) \sum_t \nabla_{\Phi} \log P(\mathbf{a}^{(t)} | \mathbf{s}^{(t)}; \Phi) b \\ &= \sum_{\tau} P(\tau; \Phi) \nabla_{\Phi} \log P(\tau; \Phi) b \\ &= \sum_{\tau} P(\tau; \Phi) \frac{\nabla_{\Phi} P(\tau; \Phi)}{P(\tau; \Phi)} b \\ &= \nabla_{\Phi} \sum_{\tau} P(\tau; \Phi) b = \nabla_{\Phi} b = 0 \end{aligned}$$

- Of what value?

Baseline II

$$\nabla_{\Phi} J(\Phi) \propto \sum_{i,t} \nabla_{\Phi} \log P(\mathbf{a}^{(i,t)} | \mathbf{s}^{(i,t)}; \Phi) \left(\sum_{t'=t}^{H^{(i)}} \gamma^{t'} R^{(i,t')} - b \right)$$

- The larger the b better, but $\sum_{t'=t}^{H^{(i)}} \gamma^{t'} R^{(i,t')} - b$ still needs to guide g_{π} to output good τ
- $\sum_{t'=t}^{H^{(i)}} \gamma^{t'} R^{(i,t')}$ an estimate of $Q_{\pi}(\mathbf{s}^{(i,t)}, \mathbf{a}^{(i,t)})$

Baseline II

$$\nabla_{\Phi} J(\Phi) \propto \sum_{i,t} \nabla_{\Phi} \log P(\mathbf{a}^{(i,t)} | \mathbf{s}^{(i,t)}; \Phi) \left(\sum_{t'=t}^{H^{(i)}} \gamma^{t'} R^{(i,t')} - b \right)$$

- The larger the b better, but $\sum_{t'=t}^{H^{(i)}} \gamma^{t'} R^{(i,t')} - b$ still needs to guide g_{π} to output good τ
- $\sum_{t'=t}^{H^{(i)}} \gamma^{t'} R^{(i,t')}$ an estimate of $Q_{\pi}(\mathbf{s}^{(i,t)}, \mathbf{a}^{(i,t)})$
- b an estimate of $V_{\pi}(\mathbf{s}^{(i,t)}) = \mathbb{E}_{\{\mathbf{s}^{(t')}, \mathbf{a}^{(t')}\}_{t'}} \left[\sum_{t'} \gamma^{t'} R^{(t')} | \mathbf{s}^{(0)} = \mathbf{s}^{(i,t)} \right]$ [3]
 - $\sum_{t'=t}^{H^{(i)}} \gamma^{t'} R^{(i,t')} - b$ estimates $Q_{\pi}(\mathbf{s}^{(i,t)}, \mathbf{a}^{(i,t)}) - V_{\pi}(\mathbf{s}^{(i,t)})$, the **advantage** of π at state $\mathbf{s}^{(i,t)}$
- In REINFORCE: $b = \frac{1}{|\{j, t'' : \mathbf{s}^{(j,t'')} = \mathbf{s}^{(i,t)}\}|} \sum_{j, t'' : \mathbf{s}^{(j,t'')} = \mathbf{s}^{(i,t)}} \sum_{t'=t''}^{H^{(j)}} \gamma^{t'} R^{(j,t')}$

Function Approximations

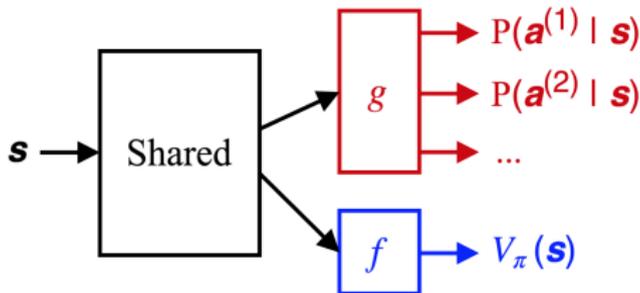
$$\nabla_{\Phi} J(\Phi) \propto \sum_{i,t} \nabla_{\Phi} \log P(\mathbf{a}^{(i,t)} | \mathbf{s}^{(i,t)}; \Phi) \left(\sum_{t'=t}^{H^{(i)}} \gamma^{t'} R^{(i,t')} - b \right)$$

- $\sum_{t'=t}^{H^{(i)}} \gamma^{t'} R^{(i,t')}$ estimates $Q_{\pi}(\mathbf{s}^{(i,t)}, \mathbf{a}^{(i,t)})$ using rolled-out from single episode
- **Actor-critic**: why not use a DNN $f_{Q_{\pi}}(\mathbf{s}, \mathbf{a}; \Theta)$ to approximate $Q_{\pi}(\mathbf{s}, \mathbf{a}), \forall \mathbf{s}, \mathbf{a}$?

Function Approximations

$$\nabla_{\Phi} J(\Phi) \propto \sum_{i,t} \nabla_{\Phi} \log P(\mathbf{a}^{(i,t)} | \mathbf{s}^{(i,t)}; \Phi) \left(\sum_{t'=t}^{H^{(i)}} \gamma^{t'} R^{(i,t')} - b \right)$$

- $\sum_{t'=t}^{H^{(i)}} \gamma^{t'} R^{(i,t')}$ estimates $Q_{\pi}(\mathbf{s}^{(i,t)}, \mathbf{a}^{(i,t)})$ using rolled-out from single episode
- **Actor-critic**: why not use a DNN $f_{Q_{\pi}}(\mathbf{s}, \mathbf{a}; \Theta)$ to approximate $Q_{\pi}(\mathbf{s}, \mathbf{a}), \forall \mathbf{s}, \mathbf{a}$?
- Baseline $b = \sum_{j: \mathbf{s}^{(j,t)} = \mathbf{s}^{(i,t)}} \sum_{t'=t}^{H^{(j)}} \gamma^{t'} R^{(j,t')}$ estimates $V_{\pi}(\mathbf{s}^{(i,t)})$
- **Advantage actor-critic**: approximates $V_{\pi}(\mathbf{s})$ with $f_{V_{\pi}}(\mathbf{s}; \Theta)$



Advantage Actor-Critic ($b = f_{V_\pi}(s; \Theta)$)

$$\nabla_{\Phi} J(\Phi) \propto \sum_{i,t} \nabla_{\Phi} \log P(\mathbf{a}^{(i,t)} | \mathbf{s}^{(i,t)}; \Phi) \left(\sum_{t'=t}^{H^{(i)}} \gamma^{t'} R^{(i,t')} - b \right)$$

- $\sum_{t'=t}^{H^{(i)}} \gamma^{t'} R^{(i,t')} \approx Q_{\pi}(\mathbf{s}^{(i,t)}, \mathbf{a}^{(i,t)})$ can be approximated by $R^{(i,t)} + \gamma f_{V_{\pi}}(\mathbf{s}^{(i,t+1)}; \Theta)$
 - No need for $f_{Q_{\pi}}$
- Bellman expectation equation for stochastic π :

$$V_{\pi}(\mathbf{s}) = \sum_{\mathbf{a}} \pi(\mathbf{a} | \mathbf{s}) \sum_{\mathbf{s}'} P(\mathbf{s}' | \mathbf{s}; \mathbf{a}) [R(\mathbf{s}, \mathbf{a}, \mathbf{s}') + \gamma V_{\pi}(\mathbf{s}')], \forall \mathbf{s}$$

- Algorithm (TD): initialize Θ and Φ arbitrarily, iterate until converge:
 - ① Take an action \mathbf{a} from \mathbf{s} using $g(\mathbf{s}; \Phi)$
 - ② Observe \mathbf{s}' and reward R , compute $\hat{Q}_{\pi} \leftarrow R + \gamma f_{V_{\pi}}(\mathbf{s}'; \Theta)$
 - ③ Update $f_{V_{\pi}}$:

$$\Theta \leftarrow \Theta - \eta \nabla_{\Theta} [\hat{Q}_{\pi} - f_{V_{\pi}}(\mathbf{s}; \Theta)]^2$$

- ④ Update g_{π} :

$$\Phi \leftarrow \Phi + \lambda \nabla_{\Phi} \log P(\mathbf{a} | \mathbf{s}; \Phi) (\hat{Q}_{\pi} - f_{V_{\pi}}(\mathbf{s}; \Theta))$$

Pitfall: Exploration

- To learn f_{V_π} based on value iteration, the agent has to *explore enough*
- But g_π is optimized for exploitation only:

$$\Phi \leftarrow \Phi + \lambda \nabla_{\Phi} \log P(\mathbf{a}|\mathbf{s}; \Phi) (\hat{Q}_\pi - f_{V_\pi}(\mathbf{s}; \Theta))$$

- Solution?

Pitfall: Exploration

- To learn f_{V_π} based on value iteration, the agent has to *explore enough*
- But g_π is optimized for exploitation only:

$$\Phi \leftarrow \Phi + \lambda \nabla_{\Phi} \log P(\mathbf{a} | \mathbf{s}; \Phi) (\hat{Q}_\pi - f_{V_\pi}(\mathbf{s}; \Theta))$$

- Solution? To maximize the entropy of $g_\pi(\mathbf{s}; \Phi)$ as well

$$\Phi \leftarrow \Phi + \lambda \nabla_{\Phi} [\log P(\mathbf{a} | \mathbf{s}; \Phi) (\hat{Q}_\pi - f_{V_\pi}(\mathbf{s}; \Theta)) + \mu H(\mathbf{a} \sim g_\pi(\mathbf{s}; \Phi))]$$

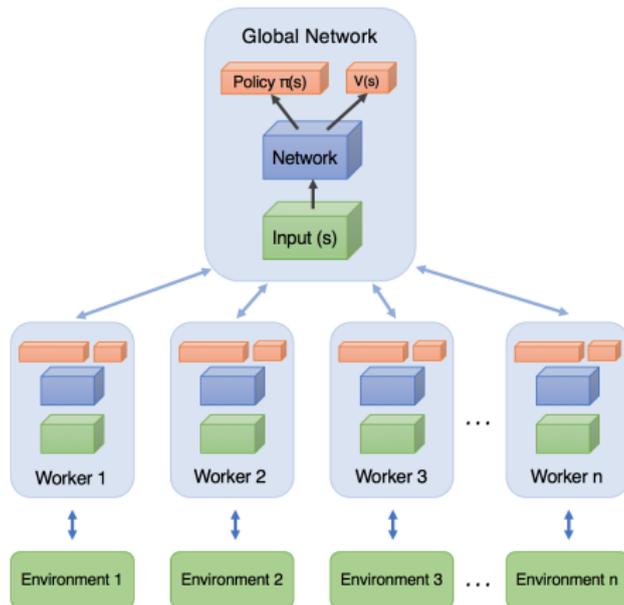
- The larger μ , the more exploration

Asynchronous Advantage Actor-Critic (A3C)

- TD estimate reduces variance at the cost of bias/divergence

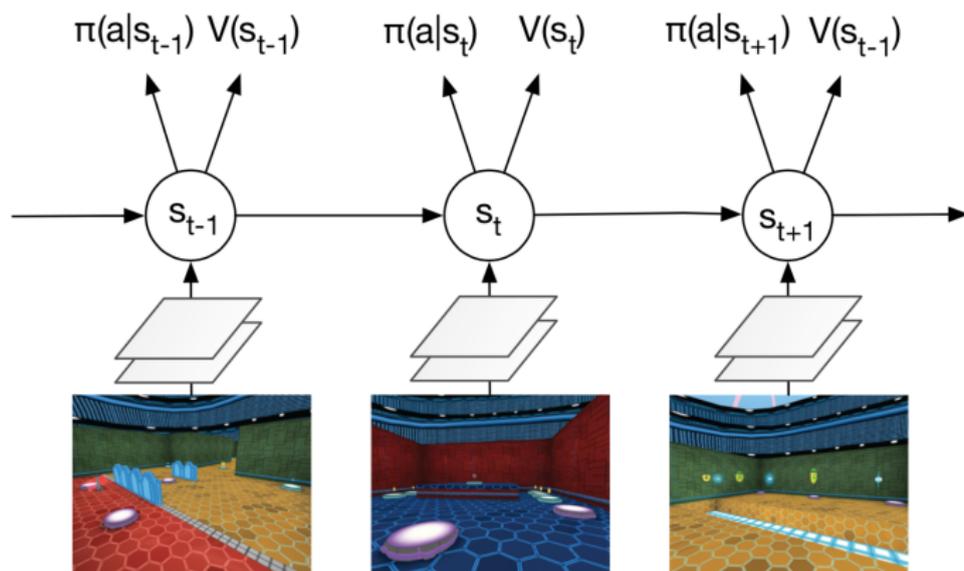
Asynchronous Advantage Actor-Critic (A3C)

- TD estimate reduces variance at the cost of bias/divergence
- A3C: use asynchronous workers to stabilize f_{V_π} training
 - An alternative to experience replay



A3C on Labyrinth

- Task: to collect apples (+1 reward) and escape (+10 reward)
- End-to-end learning from pixels to policy
- State $s^{(t)}$ modeled as a recurrent neural network (LSTM)
 - To have long-term memory



Variance Reduction by Having More Samples

Variance Reduction by Having More Samples

- Update rule for g_π in A3C:

$$\Phi \leftarrow \Phi + \lambda \nabla_{\Phi} \log P(\mathbf{a}^{(t)} | \mathbf{s}^{(t)}; \Phi) \hat{A}_\pi^{(t)}$$

where $\hat{A}_\pi^{(t)} = \hat{Q}_\pi^{(t)} - f_{V_\pi}(\mathbf{s}^{(t)}; \Theta) = (R^{(t)} + \gamma f_{V_\pi}(\mathbf{s}^{(t)}; \Theta)) - f_{V_\pi}(\mathbf{s}^{(t)}; \Theta)$

- Bellman expectation equation holds for multiple time differences:

$$\begin{aligned} Q_\pi(\mathbf{s}^{(t)}, \mathbf{a}^{(t)}) &= \mathbb{E}[R^{(t)} + \gamma V_\pi(\mathbf{s}^{(t+1)}) | \mathbf{s}^{(t)} = \mathbf{s}^{(t)}, \mathbf{a}^{(t)} = \mathbf{a}^{(t)}] \\ &= \mathbb{E}[R^{(t)} + \gamma R^{(t+1)} + \gamma^2 V_\pi(\mathbf{s}^{(t+2)})] \\ &= \mathbb{E}[R^{(t)} + \gamma R^{(t+1)} + \gamma^2 R^{(t+2)} + \gamma^3 V_\pi(\mathbf{s}^{(t+3)})] \\ &\dots \\ &= \mathbb{E}[R^{(t)} + \gamma R^{(t+1)} + \gamma^2 R^{(t+2)} + \gamma^3 R^{(t+3)} + \dots] \end{aligned}$$

- A3C replaces $\hat{A}_\pi^{(t)}$ with a K -step lookahead:

$$\hat{A}_K^{(t)} \leftarrow R^{(t)} + \gamma R^{(t+1)} + \dots + \gamma^{K-1} R^{(t+K-1)} + \gamma^K f_{V_\pi}(\mathbf{s}^{(t+K)}; \Theta) - f_{V_\pi}(\mathbf{s}^{(t)}; \Theta)$$

Variance Reduction by Having More Samples

- Update rule for g_π in A3C:

$$\Phi \leftarrow \Phi + \lambda \nabla_{\Phi} \log P(\mathbf{a}^{(t)} | \mathbf{s}^{(t)}; \Phi) \hat{A}_\pi^{(t)}$$

where $\hat{A}_\pi^{(t)} = \hat{Q}_\pi^{(t)} - f_{V_\pi}(\mathbf{s}^{(t)}; \Theta) = (R^{(t)} + \gamma f_{V_\pi}(\mathbf{s}^{(t)}; \Theta)) - f_{V_\pi}(\mathbf{s}^{(t)}; \Theta)$

- Bellman expectation equation holds for multiple time differences:

$$\begin{aligned} Q_\pi(\mathbf{s}^{(t)}, \mathbf{a}^{(t)}) &= \mathbb{E}[R^{(t)} + \gamma V_\pi(\mathbf{s}^{(t+1)}) | \mathbf{s}^{(t)} = \mathbf{s}^{(t)}, \mathbf{a}^{(t)} = \mathbf{a}^{(t)}] \\ &= \mathbb{E}[R^{(t)} + \gamma R^{(t+1)} + \gamma^2 V_\pi(\mathbf{s}^{(t+2)})] \\ &= \mathbb{E}[R^{(t)} + \gamma R^{(t+1)} + \gamma^2 R^{(t+2)} + \gamma^3 V_\pi(\mathbf{s}^{(t+3)})] \\ &\dots \\ &= \mathbb{E}[R^{(t)} + \gamma R^{(t+1)} + \gamma^2 R^{(t+2)} + \gamma^3 R^{(t+3)} + \dots] \end{aligned}$$

- A3C replaces $\hat{A}_\pi^{(t)}$ with a K -step lookahead:

$$\hat{A}_K^{(t)} \leftarrow R^{(t)} + \gamma R^{(t+1)} + \dots + \gamma^{K-1} R^{(t+K-1)} + \gamma^K f_{V_\pi}(\mathbf{s}^{(t+K)}; \Theta) - f_{V_\pi}(\mathbf{s}^{(t)}; \Theta)$$

- Update of Φ lags K time steps behind the current action

Generalized Advantage Estimation (GAE)

$$\hat{A}_K^{(t)} = R^{(t)} + \gamma R^{(t+1)} + \dots + \gamma^{K-1} R^{(t+K-1)} + \gamma^K f_{V_\pi}(\mathbf{s}^{(t+K)}; \Theta) - f_{V_\pi}(\mathbf{s}^{(t)}; \Theta)$$

- Define TD error at time t : $\delta^{(t)} = R^{(t)} + \gamma f_{V_\pi}(\mathbf{s}^{(t+1)}; \Theta) - f_{V_\pi}(\mathbf{s}^{(t)}; \Theta)$
- We have $\hat{A}_K^{(t)} = \delta^{(t)} + \gamma \delta^{(t+1)} + \dots + \gamma^{K-1} \delta^{(t+K-1)}$
- GAE [8]: let $\hat{A}_\pi^{(t)}$ be the exponential moving average of $\hat{A}_1^{(t)}, \hat{A}_2^{(t)}, \dots$:

$$\begin{aligned} \hat{A}_\pi^{(t)} &\leftarrow \hat{A}_1^{(t)} + \lambda \hat{A}_2^{(t)} + \lambda^2 \hat{A}_3^{(t)} + \dots \\ &= \delta^{(t)} + \lambda (\delta^{(t)} + \gamma \delta^{(t+1)}) + \lambda^2 (\delta^{(t)} + \gamma \delta^{(t+1)} + \gamma^2 \delta^{(t+2)}) + \dots \\ &= \frac{1}{1-\lambda} \delta^{(t)} + \frac{\lambda \gamma}{1-\lambda} \delta^{(t+1)} + \frac{\lambda^2 \gamma^2}{1-\lambda} \delta^{(t+2)} + \dots \\ &\propto \delta^{(t)} + \lambda \gamma \delta^{(t+1)} + (\lambda \gamma)^2 \delta^{(t+2)} + \dots \end{aligned}$$

Generalized Advantage Estimation (GAE)

$$\hat{A}_K^{(t)} = R^{(t)} + \gamma R^{(t+1)} + \dots + \gamma^{K-1} R^{(t+K-1)} + \gamma^K f_{V_\pi}(s^{(t+K)}; \Theta) - f_{V_\pi}(s^{(t)}; \Theta)$$

- Define TD error at time t : $\delta^{(t)} = R^{(t)} + \gamma f_{V_\pi}(s^{(t+1)}; \Theta) - f_{V_\pi}(s^{(t)}; \Theta)$
- We have $\hat{A}_K^{(t)} = \delta^{(t)} + \gamma \delta^{(t+1)} + \dots + \gamma^{K-1} \delta^{(t+K-1)}$
- GAE [8]: let $\hat{A}_\pi^{(t)}$ be the exponential moving average of $\hat{A}_1^{(t)}, \hat{A}_2^{(t)}, \dots$:

$$\begin{aligned}\hat{A}_\pi^{(t)} &\leftarrow \hat{A}_1^{(t)} + \lambda \hat{A}_2^{(t)} + \lambda^2 \hat{A}_3^{(t)} + \dots \\ &= \delta^{(t)} + \lambda (\delta^{(t)} + \gamma \delta^{(t+1)}) + \lambda^2 (\delta^{(t)} + \gamma \delta^{(t+1)} + \gamma^2 \delta^{(t+2)}) + \dots \\ &= \frac{1}{1-\lambda} \delta^{(t)} + \frac{\lambda \gamma}{1-\lambda} \delta^{(t+1)} + \frac{\lambda^2 \gamma^2}{1-\lambda} \delta^{(t+2)} + \dots \\ &\propto \delta^{(t)} + \lambda \gamma \delta^{(t+1)} + (\lambda \gamma)^2 \delta^{(t+2)} + \dots\end{aligned}$$

- Biased, but with much lower variance
 - $\delta^{(t)}$'s can have lower magnitudes when f_{V_π} is good enough
- In TD: $\hat{A}_\pi^{(t)} \leftarrow \delta^{(t)} + \lambda \gamma \delta^{(t+1)} + \dots + (\lambda \gamma)^K \delta^{(t+K)}$

Policy Optimization vs. Policy/Value Iteration

- Policy optimization:

- Policy/value-iteration-based methods (e.g., DQN, DDPG):

Policy Optimization vs. Policy/Value Iteration

- Policy optimization:
 - Optimize policy “directly”

- Policy/value-iteration-based methods (e.g., DQN, DDPG):
 - Optimize policy “indirectly” (via Q/V exploiting Bellman equations)

Policy Optimization vs. Policy/Value Iteration

- Policy optimization:
 - Optimize policy “directly”
 - More compatible with auxiliary objectives & rich NN architectures (e.g., RNN)
- Policy/value-iteration-based methods (e.g., DQN, DDPG):
 - Optimize policy “indirectly” (via Q/V exploiting Bellman equations)
 - More compatible with different exploration strategies

Policy Optimization vs. Policy/Value Iteration

- Policy optimization:
 - Optimize policy “directly”
 - More compatible with auxiliary objectives & rich NN architectures (e.g., RNN)
 - More likely to work with different tasks/settings
- Policy/value-iteration-based methods (e.g., DQN, DDPG):
 - Optimize policy “indirectly” (via Q/V exploiting Bellman equations)
 - More compatible with different exploration strategies
 - Sensitive to task/settings; but more sample-efficient when working

Reference I

- [1] Vivek F Farias and Benjamin Van Roy.
Tetris: A study of randomized constraint sampling.
In Probabilistic and Randomized Methods for Design Under Uncertainty, pages 189–201. Springer, 2006.
- [2] Meire Fortunato, Mohammad Gheshlaghi Azar, Bilal Piot, Jacob Menick, Ian Osband, Alex Graves, Vlad Mnih, Remi Munos, Demis Hassabis, Olivier Pietquin, et al.
Noisy networks for exploration.
arXiv preprint arXiv:1706.10295, 2017.
- [3] Evan Greensmith, Peter L Bartlett, and Jonathan Baxter.
Variance reduction techniques for gradient estimates in reinforcement learning.
Journal of Machine Learning Research, 5(Nov):1471–1530, 2004.

Reference II

- [4] Timothy P Lillicrap, Jonathan J Hunt, Alexander Pritzel, Nicolas Heess, Tom Erez, Yuval Tassa, David Silver, and Daan Wierstra. Continuous control with deep reinforcement learning. *arXiv preprint arXiv:1509.02971*, 2015.
- [5] Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Andrei A Rusu, Joel Veness, Marc G Bellemare, Alex Graves, Martin Riedmiller, Andreas K Fidjeland, Georg Ostrovski, et al. Human-level control through deep reinforcement learning. *Nature*, 518(7540):529–533, 2015.
- [6] Arun Nair, Praveen Srinivasan, Sam Blackwell, Cagdas Alcicek, Rory Fearon, Alessandro De Maria, Vedavyas Panneershelvam, Mustafa Suleyman, Charles Beattie, Stig Petersen, et al. Massively parallel methods for deep reinforcement learning. *arXiv preprint arXiv:1507.04296*, 2015.

Reference III

- [7] Tom Schaul, John Quan, Ioannis Antonoglou, and David Silver.
Prioritized experience replay.
arXiv preprint arXiv:1511.05952, 2015.
- [8] John Schulman, Philipp Moritz, Sergey Levine, Michael Jordan, and Pieter Abbeel.
High-dimensional continuous control using generalized advantage estimation.
arXiv preprint arXiv:1506.02438, 2015.
- [9] Hado Van Hasselt, Arthur Guez, and David Silver.
Deep reinforcement learning with double q-learning.
In *AAAI*, pages 2094–2100, 2016.
- [10] Ziyu Wang, Tom Schaul, Matteo Hessel, Hado Van Hasselt, Marc Lanctot, and Nando De Freitas.
Dueling network architectures for deep reinforcement learning.
arXiv preprint arXiv:1511.06581, 2015.

Reference IV