

# Winter School 2024

# Reinforcement Learning

## *Advanced Reinforcement Learning*

Dr. Lorenzo Brigato

Artificial Intelligence in Health and Nutrition (AIHN) Laboratory  
ARTORG Center for Biomedical Engineering Research  
University of Bern

# Outlook

- Value Function Approximation
- Policy Gradient
- Deep Reinforcement Learning

# Large-Scale Reinforcement Learning

- Reinforcement learning can be used to solve large problems, e.g.
  - Backgammon:  $10^{20}$  states
  - Go:  $10^{170}$
  - Robots: continuous state space
- How can we scale up e.g., the model-free methods for prediction and control seen before?

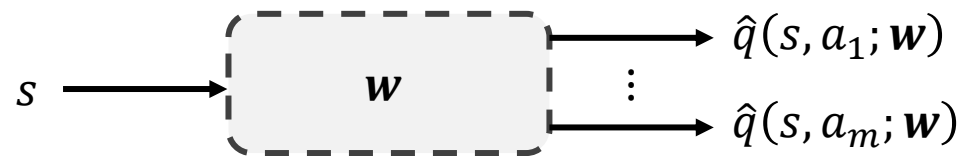
# Value Function Approximation

- So far, we have represented value function by a lookup table
  - Every state  $s$  has an entry  $V(s)$
  - Or every state-action pair  $s, a$  has an entry  $Q(s, a)$
- Problem with large MDPs:
  - There are too many states and/or actions to store in memory
  - It is too slow to learn the value of each state individually
- Solution for large MDPs:
  - Estimate value function with function approximation

$$\hat{v}(s; \mathbf{w}) \approx v_{\pi}(s)$$
$$\hat{q}(s, a; \mathbf{w}) \approx q_{\pi}(s, a)$$

- Generalize from seen states to unseen states
- Update parameter  $\mathbf{w}$  using MC or TD learning

# Types of Value Function Approximation



# Value Function Approx. By Stochastic Gradient Descent

- **Goal:** find parameter vector  $\mathbf{w}$  minimizing mean-squared error between approximate value function  $\hat{v}(s; \mathbf{w})$  and true value function  $v_\pi(s)$

$$J(\mathbf{w}) = \mathbb{E}_\pi[(v_\pi(s) - \hat{v}(s; \mathbf{w}))^2]$$

- Gradient descent finds a local minimum

$$\begin{aligned}\Delta \mathbf{w} &= -\frac{1}{2} a \nabla_{\mathbf{w}} J(\mathbf{w}) \\ &= \alpha \mathbb{E}_\pi[(v_\pi(s) - \hat{v}(s; \mathbf{w})) \nabla_{\mathbf{w}} \hat{v}(s; \mathbf{w})]\end{aligned}$$

- Stochastic gradient descent samples the gradient

$$\Delta \mathbf{w} = \alpha (v_\pi(s) - \hat{v}(s; \mathbf{w})) \nabla_{\mathbf{w}} \hat{v}(s; \mathbf{w})$$

- Expected update is equal to full gradient update

# Feature Vectors

- Represent state by a feature vector

$$x(s) = \begin{pmatrix} x_1(s) \\ \vdots \\ x_n(s) \end{pmatrix}$$

- E.g., Polynomials, Fourier Basis
- For example:
  - Distance of robot from landmarks
  - Trends in the stock market
  - Piece and pawn configurations in chess



# Linear Value Function Approximation

- Represent value function by a linear combination of features

$$\hat{v}(s; \mathbf{w}) = x(s)^T \mathbf{w} = \sum_{i=0}^n x_i(s) w_i$$

- Objective function is quadratic in parameters  $\mathbf{w}$

$$J(\mathbf{w}) = \mathbb{E}_{\pi}[(v_{\pi}(s) - x(s)^T \mathbf{w})^2]$$

- Stochastic gradient descent converges on global optimum
- Update rule is particularly simple

$$\begin{aligned} \nabla_{\mathbf{w}} \hat{v}(s; \mathbf{w}) &= x(s) \\ \Delta \mathbf{w} &= \alpha (v_{\pi}(s) - \hat{v}(s; \mathbf{w})) x(s) \end{aligned}$$

- Update = step-size  $\times$  prediction error  $\times$  feature value



# Incremental Prediction Algorithms

- Assumed true value  $v_{\pi}(s)$  function given by supervisor
- But in RL there is no supervisor, only rewards
- In practice, we substitute a target for  $v_{\pi}(s)$ 
  - For MC, the target is the return  $G_t$

$$\Delta \mathbf{w} = \alpha (G_t - \hat{v}(S_t; \mathbf{w})) \nabla_{\mathbf{w}} \hat{v}(S_t; \mathbf{w})$$

- For TD(0), the target is the TD target  $G_{t:t+1} = R_{t+1} + \gamma \hat{v}(S_{t+1}; \mathbf{w})$

$$\Delta \mathbf{w} = \alpha (R_{t+1} + \gamma \hat{v}(S_{t+1}; \mathbf{w}) - \hat{v}(S_t; \mathbf{w})) \nabla_{\mathbf{w}} \hat{v}(S_t; \mathbf{w})$$

# Monte-Carlo with Value Function Approximation (1/2)

- Return  $G_t$  is an unbiased, noisy sample of true value  $v_\pi(s_t)$
- Can therefore apply supervised learning to “training data”:

$$\langle S_1, G_1 \rangle, \langle S_2, G_2 \rangle, \dots, \langle S_T, G_T \rangle$$

- For example, using linear Monte-Carlo policy evaluation

$$\begin{aligned}\Delta \mathbf{w} &= \alpha(G_t - \hat{v}(S_t; \mathbf{w})) \nabla_{\mathbf{w}} \hat{v}(S_t; \mathbf{w}) \\ &= \alpha(G_t - \hat{v}(S_t; \mathbf{w})) \mathbf{x}(S_t)\end{aligned}$$

- Monte-Carlo evaluation converges to a local optimum
- Even when using non-linear value function approximation

# Monte-Carlo with Value Function Approximation (2/2)

## Gradient Monte Carlo Algorithm for Estimating $\hat{v} \approx v_\pi$

Input: the policy  $\pi$  to be evaluated

Input: a differentiable function  $\hat{v} : \mathcal{S} \times \mathbb{R}^d \rightarrow \mathbb{R}$

Algorithm parameter: step size  $\alpha > 0$

Initialize value-function weights  $\mathbf{w} \in \mathbb{R}^d$  arbitrarily (e.g.,  $\mathbf{w} = \mathbf{0}$ )

Loop forever (for each episode):

    Generate an episode  $S_0, A_0, R_1, S_1, A_1, \dots, R_T, S_T$  using  $\pi$

    Loop for each step of episode,  $t = 0, 1, \dots, T - 1$ :

$$\mathbf{w} \leftarrow \mathbf{w} + \alpha [G_t - \hat{v}(S_t, \mathbf{w})] \nabla \hat{v}(S_t, \mathbf{w})$$

# TD Learning with Value Function Approximation (1/2)

- The TD-target  $G_{t:t+1} = R_{t+1} + \gamma \hat{v}(s_{t+1}; \mathbf{w})$  is a biased sample of true value  $v_{\pi}(s_t)$
- Can still apply supervised learning to “training data”:

$$\langle S_1, R_1 + \gamma \hat{v}(S_2; \mathbf{w}) \rangle, \langle S_2, R_3 + \gamma \hat{v}(S_3; \mathbf{w}) \rangle, \dots, \langle S_{T-1}, R_T \rangle$$

- For example, using linear TD(0)

$$\begin{aligned} \Delta \mathbf{w} &= \alpha (R + \gamma \hat{v}(S'; \mathbf{w}) - \hat{v}(S_t; \mathbf{w})) \nabla_{\mathbf{w}} \hat{v}(S_t; \mathbf{w}) \\ &= \alpha \delta \mathbf{x}(S) \end{aligned}$$

- (Semi-Gradient) We do not consider the effect of changing  $\mathbf{w}$  on the target
- Linear TD(0) converges (close) to global optimum

# TD Learning with Value Function Approximation (2/2)

## Semi-gradient TD(0) for estimating $\hat{v} \approx v_\pi$

Input: the policy  $\pi$  to be evaluated

Input: a differentiable function  $\hat{v} : \mathcal{S}^+ \times \mathbb{R}^d \rightarrow \mathbb{R}$  such that  $\hat{v}(\text{terminal}, \cdot) = 0$

Algorithm parameter: step size  $\alpha > 0$

Initialize value-function weights  $\mathbf{w} \in \mathbb{R}^d$  arbitrarily (e.g.,  $\mathbf{w} = \mathbf{0}$ )

Loop for each episode:

  Initialize  $S$

  Loop for each step of episode:

    Choose  $A \sim \pi(\cdot|S)$

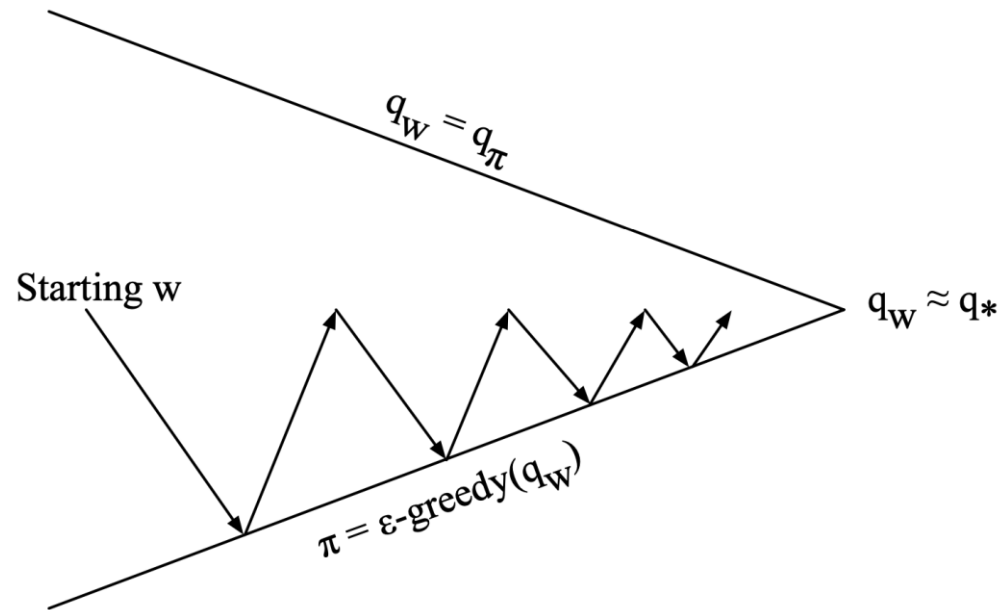
    Take action  $A$ , observe  $R, S'$

$\mathbf{w} \leftarrow \mathbf{w} + \alpha [R + \gamma \hat{v}(S', \mathbf{w}) - \hat{v}(S, \mathbf{w})] \nabla \hat{v}(S, \mathbf{w})$

$S \leftarrow S'$

  until  $S$  is terminal

# Control with Value Function Approximation



- Policy evaluation
  - Approximate policy evaluation,  $\hat{q}(\cdot, \cdot; w) \approx q_\pi$
- Policy improvement
  - $\epsilon$ -Greedy policy improvement

# Action-Value Function Approximation

- Approximate the action-value function

$$\hat{q}(S, A; \mathbf{w}) \approx q_{\pi}(S, A)$$

- Minimize mean-squared error between approximate action-value function  $\hat{q}(S, A, \mathbf{w})$  and true action-value function  $q_{\pi}(S, A)$

$$J(\mathbf{w}) = \mathbb{E}_{\pi}[(q_{\pi}(S, A) - \hat{q}(S, A; \mathbf{w}))^2]$$

- Use stochastic gradient descent to find a local minimum

$$\begin{aligned} -\frac{1}{2} \nabla_{\mathbf{w}} J(\mathbf{w}) &= (q_{\pi}(S, A) - \hat{q}(S, A; \mathbf{w})) \nabla_{\mathbf{w}} \hat{q}(S, A; \mathbf{w}) \\ \Delta \mathbf{w} &= \alpha (q_{\pi}(S, A) - \hat{q}(S, A; \mathbf{w})) \nabla_{\mathbf{w}} \hat{q}(S, A; \mathbf{w}) \end{aligned}$$

# Linear Action-Value Function Approximation

- Represent state and action by a feature vector

$$x(S, A) = \begin{pmatrix} x_1(S, A) \\ \vdots \\ x_n(S, A) \end{pmatrix}$$

- Represent action-value function by linear combination of features

$$\hat{q}(S, A, \mathbf{w}) = x(S, A)^T \mathbf{w} = \sum_{i=0}^n x_i(S, A) w_i$$

- Stochastic gradient descent update

$$\begin{aligned} \nabla_{\mathbf{w}} \hat{q}(S, A; \mathbf{w}) &= x(S, A) \\ \Delta \mathbf{w} &= \alpha (q_{\pi}(S, A) - \hat{q}(S, A; \mathbf{w})) x(S, A) \end{aligned}$$



# Incremental Control Algorithms

- Like prediction, we must substitute a target for  $q_\pi(S, A)$ 
  - For MC, the target is the return  $G_t$

$$\Delta \mathbf{w} = \alpha (G_t - \hat{q}(S_t, A_t; \mathbf{w})) \nabla_{\mathbf{w}} \hat{q}(S_t, A_t; \mathbf{w})$$

- For TD(0), the target is the TD target  $G_{t:t+1} = R_{t+1} + \gamma \hat{v}(s_{t+1}; \mathbf{w})$

$$\Delta \mathbf{w} = \alpha (R_{t+1} + \gamma \hat{v}(s_{t+1}; \mathbf{w}) - \hat{q}(S_t, A_t; \mathbf{w})) \nabla_{\mathbf{w}} \hat{q}(S_t, A_t; \mathbf{w})$$

# Policy-Based Reinforcement Learning

- Previously we approximated the value or action-value function using parameters  $w$

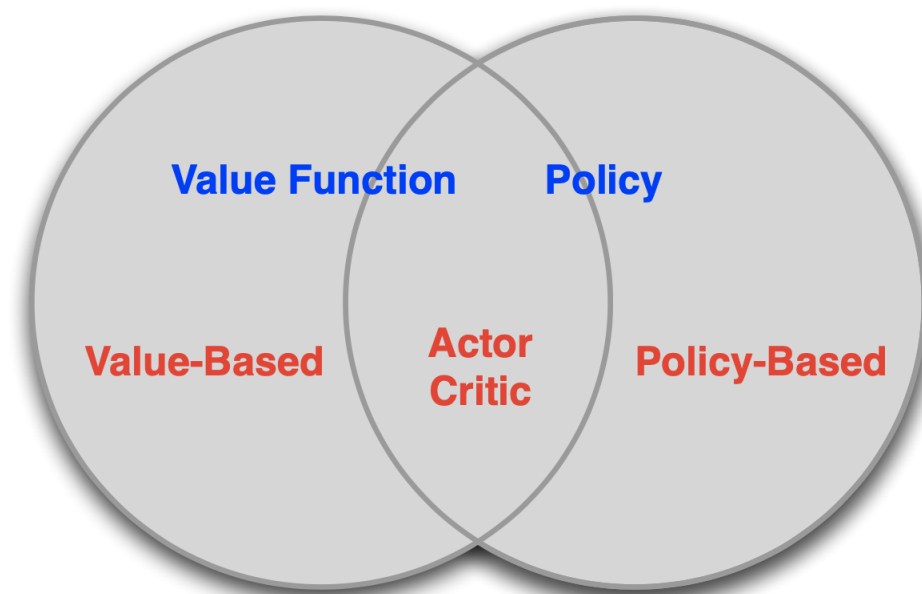
$$\begin{aligned}v_w(s) &\approx v_\pi(s) \\ q_w(s, a) &\approx q_\pi(s, a)\end{aligned}$$

- A policy was generated directly from the value function
  - e.g., using  $\epsilon$ -greedy
- We now directly parametrize the policy

$$\pi_\theta(s, a) = \mathbb{P}[a \mid s; \theta]$$

# Value-Based and Policy-Based RL

- Value Based
  - Learnt Value Function
  - Implicit policy (e.g.,  $\epsilon$ -greedy)
- Policy Based
  - No Value Function
  - Learnt Policy
- Actor-Critic
  - Learnt Value Function
  - Learnt Policy



# Advantages of Policy-Based RL

## — Advantages:

- Better convergence properties (in contrast to value function approximation that can oscillate in some configurations)
- Effective in high-dimensional or continuous action spaces
- Can learn stochastic policies

## — Disadvantages:

- Typically converge to a local rather than global optimum
- Evaluating a policy is typically inefficient and high variance

# Policy Optimization

- Policy based reinforcement learning is an **optimization** problem
- Find  $\theta$  that maximizes  $J(\theta)$
- Many possible optimization approaches
  - Gradient-free (e.g., Hill climbing, Genetic algorithms, etc.)
  - Gradient-based (e.g., Gradient descent)
- We focus on gradient descent, many extensions possible

# Policy Gradient

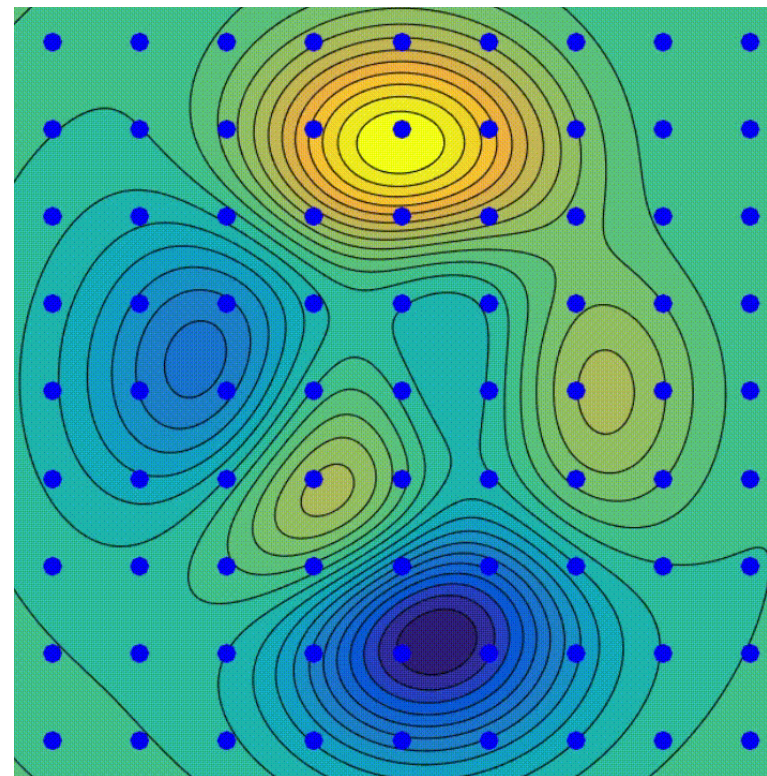
- Let  $J(\theta)$  be any policy objective function
- Policy gradient algorithms search for a local maximum in  $J(\theta)$  by **ascending the gradient** of the policy, w.r.t. parameters  $\theta$

$$\Delta\theta = a\nabla_{\theta}J(\boldsymbol{\theta})$$

- Where  $\nabla_{\theta}J(\boldsymbol{\theta})$  is the **policy gradient**

$$\nabla_{\theta}J(\boldsymbol{\theta}) = \begin{pmatrix} \frac{\partial J(\theta)}{\partial \theta_1} \\ \vdots \\ \frac{\partial J(\theta)}{\partial \theta_n} \end{pmatrix}$$

- and  $\alpha$  is a step-size parameter



# Policy Objective Functions

- Goal: given policy  $\pi_\theta(s, a)$  with parameters  $\theta$ , find best  $\theta$
- But how do we measure the quality of a policy  $\pi_\theta$ ?
- In episodic environments we can use the **start value**

$$J_1(\theta) = V^{\pi_\theta}(s_1) = \mathbb{E}_{\pi_\theta}[V_1]$$

- In continuing environments, we can use the **average value**

$$J_{avV}(\theta) = \sum_s d^{\pi_\theta}(s) V^{\pi_\theta}(s)$$

- Or the **average reward per time-step**

$$J_{avR}(\theta) = \sum_s d^{\pi_\theta}(s) \sum_a \pi_\theta(s, a) \mathcal{R}_s^a$$

where  $d^{\pi_\theta}(s)$  is **stationary distribution** of Markov chain for  $\pi_\theta$

# One-Step MDPs

- Consider a simple class of **one-step** MDPs
  - Starting in state  $s \sim d(s)$
  - Terminating after one time-step with reward  $r = \mathcal{R}_s^a$
- Use likelihood ratios to compute the policy gradient

$$\begin{aligned}\nabla_{\theta} \pi_{\theta}(s, a) &= \pi_{\theta}(s, a) \frac{\nabla_{\theta} \pi_{\theta}(s, a)}{\pi_{\theta}(s, a)} \\ &= \pi_{\theta}(s, a) \nabla_{\theta} \log \pi_{\theta}(s, a)\end{aligned}$$

$$\begin{aligned}J(\theta) &= \mathbb{E}_{\pi_{\theta}}[r] \\ &= \sum_{s \in \mathcal{S}} d(s) \sum_{a \in \mathcal{A}} \pi_{\theta}(s, a) \mathcal{R}_s^a \\ \nabla_{\theta} J(\theta) &= \sum_{s \in \mathcal{S}} d(s) \sum_{a \in \mathcal{A}} \pi_{\theta}(s, a) \nabla_{\theta} \log \pi_{\theta}(s, a) \mathcal{R}_s^a \\ &= \mathbb{E}_{\pi_{\theta}}[\nabla_{\theta} \log \pi_{\theta}(s, a) r]\end{aligned}$$



# Policy Gradient Theorem

- The policy gradient theorem generalizes the likelihood ratio approach to multi-step MDPs
- Replaces instantaneous reward  $r$  with long-term value  $q_\pi(s, a)$
- Policy gradient theorem applies to start state objective, average reward and average value objective

## Theorem (Policy Gradient)

*For any differentiable policy  $\pi_\theta(s, a)$  and for any of the policy objective functions  $J = J_1, J_{avR}$  or  $\frac{1}{1-\gamma}J_{avV}$  the policy gradient is*

$$\nabla_\theta J(\theta) = \mathbb{E}_{\pi_\theta} [\nabla_\theta \log \pi_\theta(s, a) q_\pi(s, a) ]$$

# Monte-Carlo Policy Gradient (REINFORCE)

- Update parameters by **stochastic** gradient ascent
- Using **policy gradient theorem**
- Using **return**  $G_t$  as an unbiased sample of  $q_{\pi_{\theta}}(s_t, a_t)$

$$\Delta\theta_t = \alpha \nabla_{\theta} \log \pi_{\theta}(s_t, a_t) G_t$$

## REINFORCE: Monte-Carlo Policy-Gradient Control (episodic) for $\pi_*$

Input: a differentiable policy parameterization  $\pi(a|s, \theta)$

Algorithm parameter: step size  $\alpha > 0$

Initialize policy parameter  $\theta \in \mathbb{R}^{d'}$  (e.g., to  $\mathbf{0}$ )

Loop forever (for each episode):

    Generate an episode  $S_0, A_0, R_1, \dots, S_{T-1}, A_{T-1}, R_T$ , following  $\pi(\cdot|\cdot, \theta)$

    Loop for each step of the episode  $t = 0, 1, \dots, T - 1$ :

$$G \leftarrow \sum_{k=t+1}^T \gamma^{k-t-1} R_k \quad (G_t)$$

$$\theta \leftarrow \theta + \alpha \gamma^t G \nabla \ln \pi(A_t|S_t, \theta)$$

# Reducing Variance Using a Critic

- Monte-Carlo policy gradient still has high variance
- We use a **critic** to estimate the action-value function

$$q_w(s, a) \approx q_{\pi_\theta}(s, a)$$

- Actor-critic algorithms maintain two sets of parameters
  - **Critic** Updates action-value function parameters  $w$
  - **Actor** Updates policy parameters  $\theta$ , in direction suggested by critic
- Actor-critic algorithms follow an approximate policy gradient

$$\nabla_\theta J(\theta) \approx \mathbb{E}_{\pi_\theta} [\nabla_\theta \log \pi_\theta(s, a) q_w(s, a) ]$$

$$\Delta\theta = \alpha \nabla_\theta \log \pi_\theta(s, a) q_w(s, a)$$

# Estimating the Action-Value Function

- The critic is solving a familiar problem: policy evaluation (prediction)
- How good is policy  $\pi_\theta$  for current parameters  $\theta$ ?
- Could also use e.g., least-squares policy evaluation

# Action-Value Actor-Critic (1/2)

- Simple actor-critic algorithm based on action-value critic
- Using linear value function approximation  $q_w(s, a) = \varphi(s, a)^T w$ 
  - Critic Updates  $w$  by linear TD(0)
  - Actor Updates  $\theta$  by policy gradient

# Action-Value Actor-Critic (2/2)

## One-step Actor-Critic (episodic), for estimating $\pi_{\theta} \approx \pi_*$

Input: a differentiable policy parameterization  $\pi(a|s, \theta)$

Input: a differentiable state-value function parameterization  $\hat{v}(s, \mathbf{w})$

Parameters: step sizes  $\alpha^{\theta} > 0, \alpha^{\mathbf{w}} > 0$

Initialize policy parameter  $\theta \in \mathbb{R}^{d'}$  and state-value weights  $\mathbf{w} \in \mathbb{R}^d$  (e.g., to  $\mathbf{0}$ )

Loop forever (for each episode):

  Initialize  $S$  (first state of episode)

$I \leftarrow 1$

  Loop while  $S$  is not terminal (for each time step):

$A \sim \pi(\cdot|S, \theta)$

    Take action  $A$ , observe  $S', R$

$\delta \leftarrow R + \gamma \hat{v}(S', \mathbf{w}) - \hat{v}(S, \mathbf{w})$       (if  $S'$  is terminal, then  $\hat{v}(S', \mathbf{w}) \doteq 0$ )

$\mathbf{w} \leftarrow \mathbf{w} + \alpha^{\mathbf{w}} \delta \nabla \hat{v}(S, \mathbf{w})$

$\theta \leftarrow \theta + \alpha^{\theta} I \delta \nabla \ln \pi(A|S, \theta)$

$I \leftarrow \gamma I$

$S \leftarrow S'$

# Approximation with Deep Networks

- So far, the feature representation was typically “fixed”
- The parametrised functions  $\hat{v}(s; w) / \pi_{\theta}(s, a)$  were linear mappings of input features
- More complicated non-linear mappings are needed to generalize to more complex domains
- Popular choice is to use deep neural networks
  - Known to discover useful feature representation tailored to the specific task
  - We can leverage extensive research on architectures and optimisation from Supervised Learning

# Batch Reinforcement Learning

- Gradient descent is simple and appealing
- But it is not sample efficient
- Batch methods seek to find the best fitting value function
- Given the agent's experience (“training data”)



# Stochastic Gradient Descent with Experience Replay

- Given experience consisting of  $\langle state, value \rangle$  pairs

$$\mathcal{D} = \{ \langle s_1, v_1^\pi \rangle, \langle s_2, v_2^\pi \rangle, \dots, \langle s_T, v_T^\pi \rangle \}$$

- Repeat:

1. Sample state, value from experience

$$\langle s, v^\pi \rangle \sim \mathcal{D}$$

2. Apply stochastic gradient descent update

$$\Delta \mathbf{w} = \alpha (v_\pi(s) - \hat{v}(s; \mathbf{w})) \nabla_{\mathbf{w}} \hat{v}(s; \mathbf{w})$$

- Converges to least squares solution

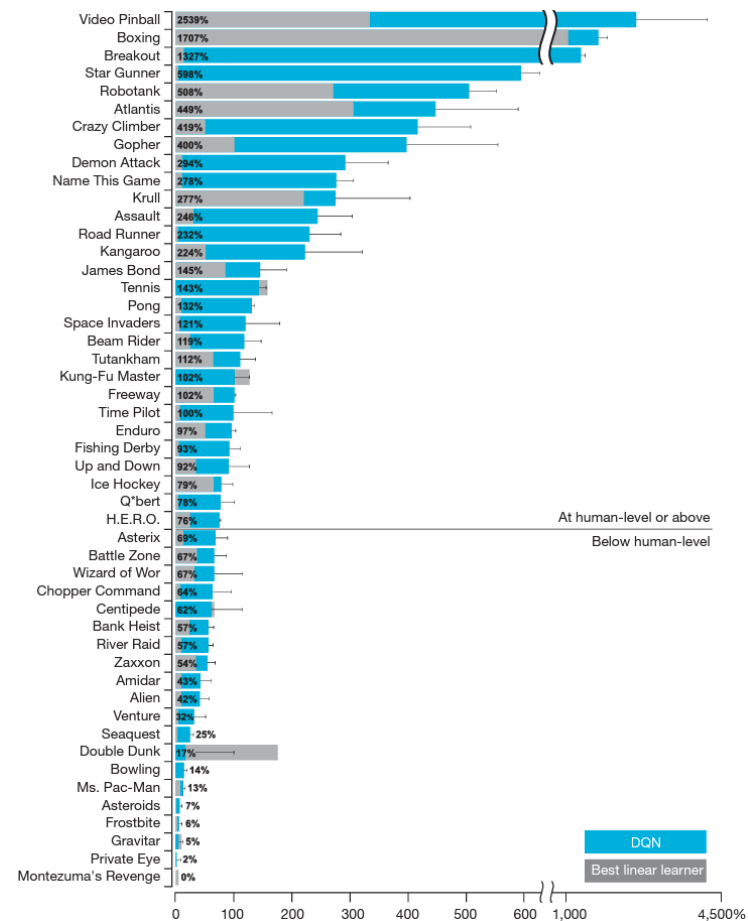
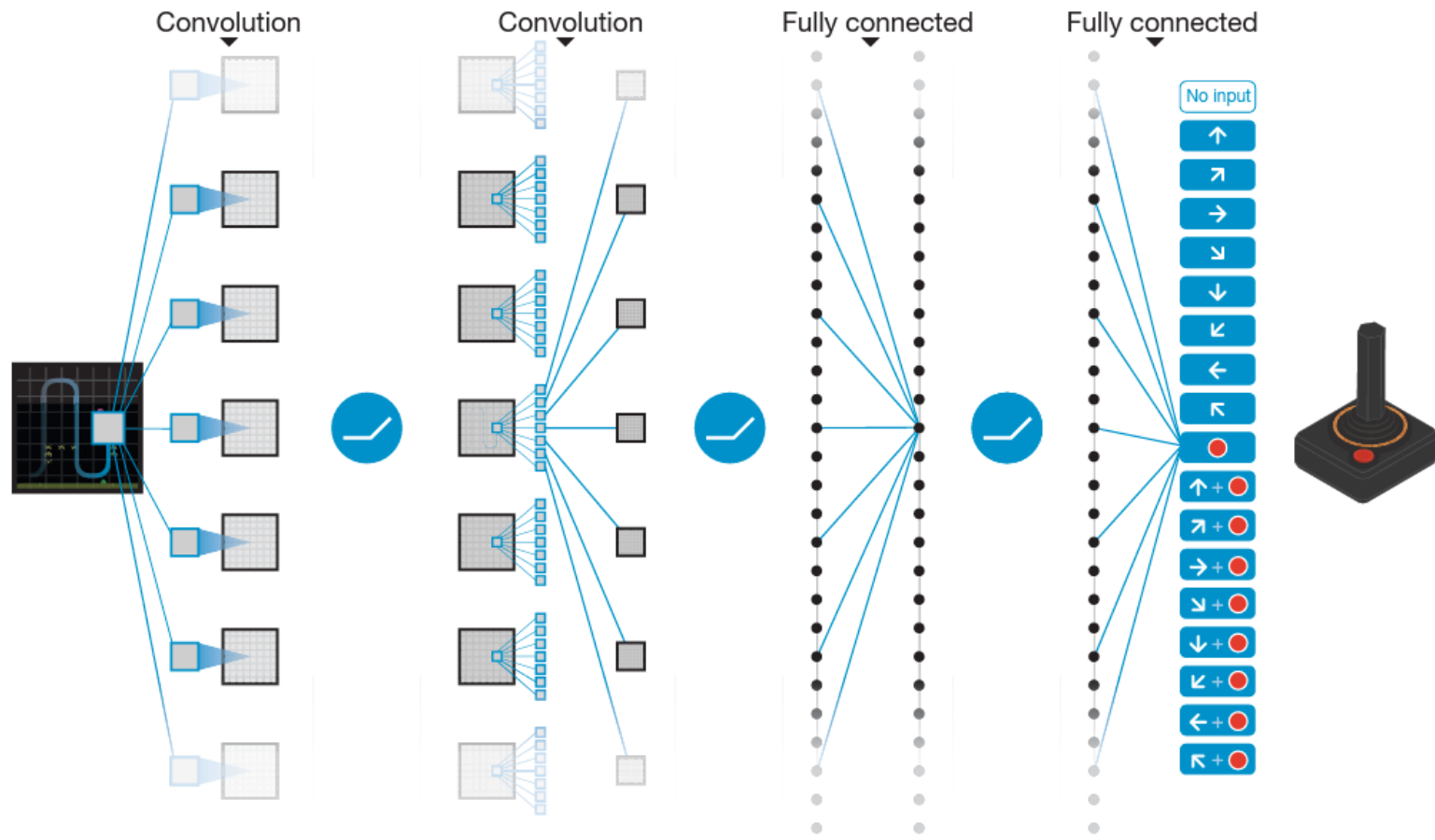
$$\mathbf{w}^\pi = \underset{\mathbf{w}}{\operatorname{argmin}} LS(\mathbf{w})$$

# Experience Replay in Deep Q-Networks (DQN)

- Example of DQN that uses experience replay and fixed Q-targets
  - Take action  $a_t$  according to  $\epsilon$ -greedy policy
  - Store transition  $(s_t, a_t, r_{t+1}, s_{t+1})$  in replay memory  $\mathcal{D}$
  - Sample random mini-batch of transitions  $(s, a, r, s')$  from  $\mathcal{D}$
  - Compute Q-learning targets w.r.t. the parameters of a DQN  $w^-$
  - Optimize MSE between Q-network and Q-learning targets

$$\mathcal{L}_i(w_i) = \mathbb{E}_{s,a,r,s' \sim \mathcal{D}_i} \left[ \left( r + \gamma \max_{a'} Q(s', a'; w_i^-) - Q(s, a; w_i) \right)^2 \right]$$

# DQN in Atari



[doi:10.1038/nature14236]

# Continuous Action Spaces

- Vanilla DQN can't be used for continuous action spaces (CAS)
  - Too many outputs to parametrize with a neural network
  - Discretization would be needed (suboptimal)
- Deep Deterministic Policy Gradient (DDPG) extends DQN to CAS
  - DDPG is an Actor Critic algorithm with critic  $Q_w$  and actor  $\pi_\theta$  networks
  - Policy is deterministic, noise added for exploration  $a_t = \pi_\theta(s_t) + \epsilon$
  - Update the critic with the actor as max and then the actor to maximize the critic

$$\mathcal{L}_i(w_i) = \mathbb{E}_{s,a,r,s' \sim \mathcal{D}_i} [(r + \gamma Q(s', \pi(s'; \theta_i^-); w_i^-) - Q(s, a; w_i))^2]$$

$$\mathcal{L}_i(\theta_i) = \mathbb{E}_{s \sim \mathcal{D}_i} [Q(s, \pi(s; \theta); w_i)]$$

# DDPG

---

**Algorithm 1** Deep Deterministic Policy Gradient

---

- 1: Input: initial policy parameters  $\theta$ , Q-function parameters  $\phi$ , empty replay buffer  $\mathcal{D}$
- 2: Set target parameters equal to main parameters  $\theta_{\text{targ}} \leftarrow \theta$ ,  $\phi_{\text{targ}} \leftarrow \phi$
- 3: **repeat**
- 4:   Observe state  $s$  and select action  $a = \text{clip}(\mu_{\theta}(s) + \epsilon, a_{\text{Low}}, a_{\text{High}})$ , where  $\epsilon \sim \mathcal{N}$
- 5:   Execute  $a$  in the environment
- 6:   Observe next state  $s'$ , reward  $r$ , and done signal  $d$  to indicate whether  $s'$  is terminal
- 7:   Store  $(s, a, r, s', d)$  in replay buffer  $\mathcal{D}$
- 8:   If  $s'$  is terminal, reset environment state.
- 9:   **if** it's time to update **then**
- 10:     **for** however many updates **do**
- 11:       Randomly sample a batch of transitions,  $B = \{(s, a, r, s', d)\}$  from  $\mathcal{D}$
- 12:       Compute targets

$$y(r, s', d) = r + \gamma(1 - d)Q_{\phi_{\text{targ}}}(s', \mu_{\theta_{\text{targ}}}(s'))$$

- 13:       Update Q-function by one step of gradient descent using

$$\nabla_{\phi} \frac{1}{|B|} \sum_{(s,a,r,s',d) \in B} (Q_{\phi}(s, a) - y(r, s', d))^2$$

- 14:       Update policy by one step of gradient ascent using

$$\nabla_{\theta} \frac{1}{|B|} \sum_{s \in B} Q_{\phi}(s, \mu_{\theta}(s))$$

- 15:       Update target networks with

$$\begin{aligned}\phi_{\text{targ}} &\leftarrow \rho\phi_{\text{targ}} + (1 - \rho)\phi \\ \theta_{\text{targ}} &\leftarrow \rho\theta_{\text{targ}} + (1 - \rho)\theta\end{aligned}$$

- 16:     **end for**
- 17:   **end if**
- 18: **until** convergence