# Actor-Critic Methods

Alina Vereshchaka

CSE4/546 Reinforcement Learning
Spring 2021

*avereshc@buffalo.edu*

April 12, 2021

*Slides are adopted from Deep Reinforcement Learning by Sergey Levine & Policy Gradients by David Silver

# Table of Contents

$$\theta^* = \arg\max_{\theta} R_{\tau \sim p_\theta(\tau)} \left[ \sum_t r(s_t, a_t) \right]$$

- Model-based RL:

$$\theta^* = \arg\max_{\theta} R_{\tau \sim p_\theta(\tau)} \left[ \sum_t r(s_t, a_t) \right]$$

- Model-based RL: estimate the transition model and then:
    - Use it for planning (no explicit policy)
    - Use it to improve a policy

$$\theta^* = \arg\max_\theta R_{\tau \sim p_\theta(\tau)} \left[ \sum_t r(s_t, a_t) \right]$$

- Model-based RL: estimate the transition model and then:
  - Use it for planning (no explicit policy)
  - Use it to improve a policy
- Value-based:

$$\theta^* = \arg\max_{\theta} R_{\tau \sim p_\theta(\tau)} \left[ \sum_t r(s_t, a_t) \right]$$

- Model-based RL: estimate the transition model and then:
  - Use it for planning (no explicit policy)
  - Use it to improve a policy
- Value-based: estimate value function or Q-function of the current policy (no explicit policy)
- Policy-gradient:

# Types of RL algorithms

$$\theta^* = \arg\max_\theta R_{\tau \sim p_\theta(\tau)} \left[ \sum_t r(s_t, a_t) \right]$$
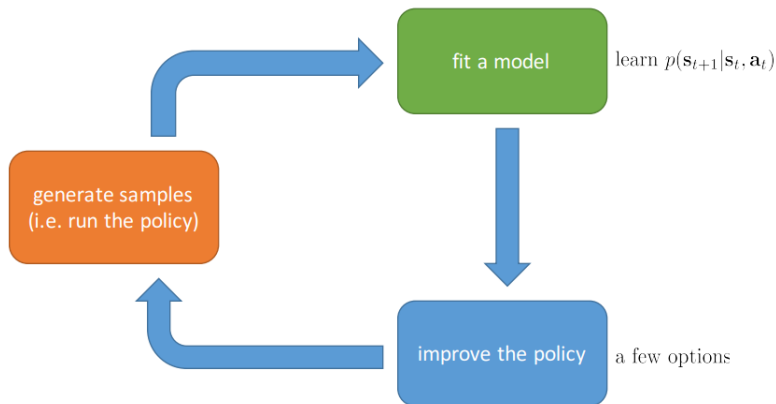
- **Model-based RL:** estimate the transition model and then:
  - Use it for planning (no explicit policy)
  - Use it to improve a policy
- **Value-based:** estimate value function or Q-function of the current policy (no explicit policy)
- **Policy-gradient:** directly differentiate the objective
- **Actor-critic:**

$$\theta^* = \arg\max_{\theta} R_{\tau \sim p_\theta(\tau)} \left[ \sum_t r(s_t, a_t) \right]$$

- Model-based RL: estimate the transition model and then:
  - Use it for planning (no explicit policy)
  - Use it to improve a policy

- Value-based: estimate value function or Q-function of the current policy (no explicit policy)

- Policy-gradient: directly differentiate the objective

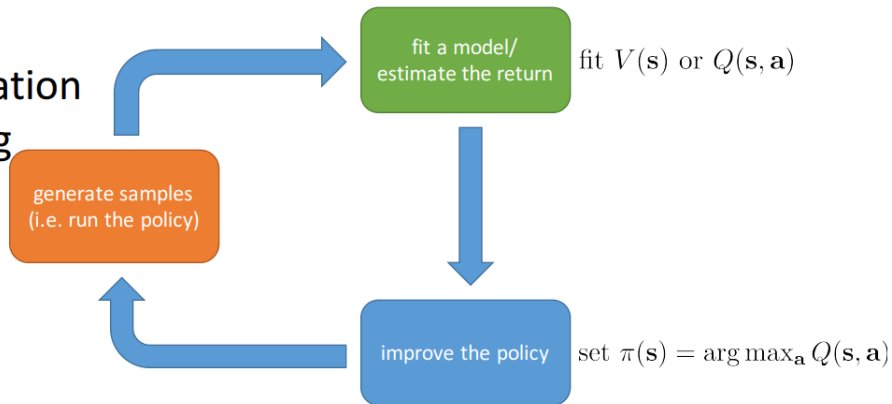- Actor-critic: estimate value function or Q-function of the current policy, use it to improve the policy

Examples:
- Value-Iteration
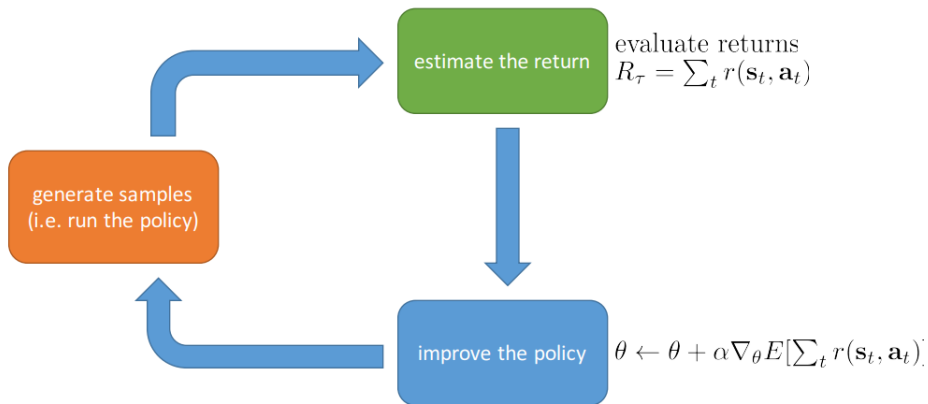- Q-Learning
- DQN



fit a model/ estimate the return — fit $V(\mathbf{s})$ or $Q(\mathbf{s}, \mathbf{a})$

generate samples (i.e. run the policy)

improve the policy — set $\pi(\mathbf{s}) = \arg\max_{\mathbf{a}} Q(\mathbf{s}, \mathbf{a})$

estimate the return

evaluate returns
$R_\tau = \sum_t r(\mathbf{s}_t, \mathbf{a}_t)$

generate samples
(i.e. run the policy)

improve the policy

$\theta \leftarrow \theta + \alpha \nabla_\theta E[\sum_t r(\mathbf{s}_t, \mathbf{a}_t)]$

fit $V(\mathbf{s})$ or $Q(\mathbf{s}, \mathbf{a})$

evaluate returns using $V$ or $Q$!

fit a model/ estimate the return

generate samples (i.e. run the policy)

improve the policy $\quad \theta \leftarrow \theta + \alpha \nabla_\theta E[\sum_t r(\mathbf{s}_t, \mathbf{a}_t)]$

- **Sample efficiency:** How many samples do we need to get a good policy?

- **Sample efficiency:** How many samples do we need to get a good policy?
- Is the algorithm off/on policy?
    - Off policy: able to improve the policy without generating new samples from that policy
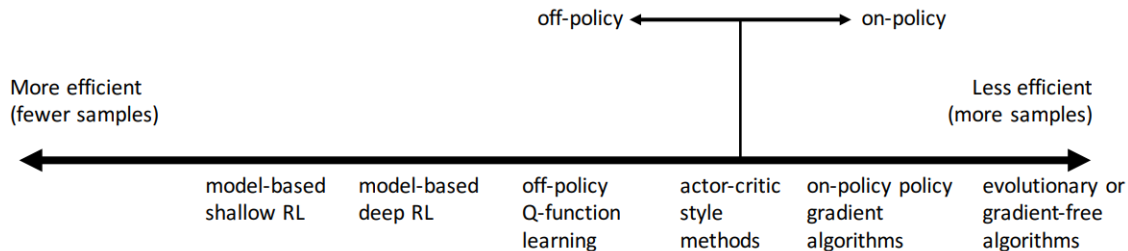
- **Sample efficiency:** How many samples do we need to get a good policy?
- Is the algorithm off/on policy?
    - Off policy: able to improve the policy without generating new samples from that policy
    - On policy: each time the policy is changed, even a little bit, we need to generate new samples

# REINFORCE (Monte-Carlo Policy Gradient)

- ‣ Update parameters by stochastic gradient ascent

- ‣ Using policy gradient theorem

- ‣ Using return $G_t$ as an unbiased sample of $Q^{\pi_\theta}(s_t, a_t)$

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

---

**REINFORCE, A Monte-Carlo Policy-Gradient Method (episodic)**

Input: a differentiable policy parameterization $\pi(a|s, \boldsymbol{\theta}), \forall a \in \mathcal{A}, s \in \mathcal{S}, \boldsymbol{\theta} \in \mathbb{R}^n$

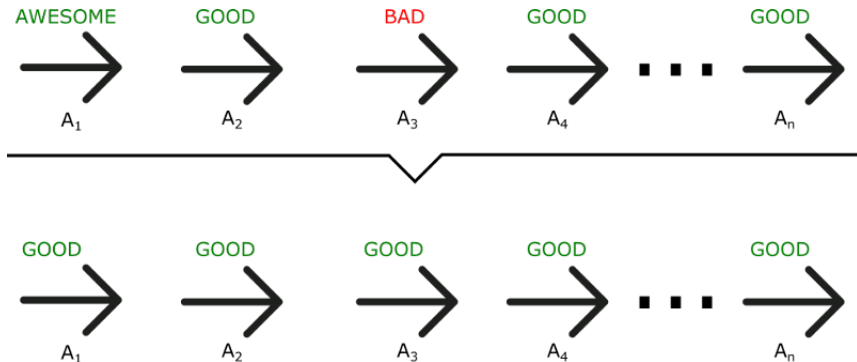Initialize policy weights $\boldsymbol{\theta}$

Repeat forever:

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

　For each step of the episode $t = 0, \ldots, T - 1$:

　　$G_t \leftarrow$ return from step $t$

　　$\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} + \alpha \gamma^t G_t \nabla_{\boldsymbol{\theta}} \log \pi(A_t|S_t, \boldsymbol{\theta})$

---

Policy Update: $\quad \Delta \theta = \alpha * \nabla_\theta * (log\ \pi(S_t, A_t, \theta)) * \cancel{R(t)}$

New update: $\Delta \theta = \alpha * \nabla_\theta * (log\ \pi(S_t, A_t, \theta)) * \boxed{Q(S_t, A_t)}$

# Table of Contents

# Actor-Critic

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

$$Q_w(s, a) \approx Q_{\pi_\theta}(s, a)$$

# Actor-Critic

- Monte-Carlo policy gradient still has high variance

- We can 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-Critic

- Monte-Carlo policy gradient still has high variance

- We can 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

- Monte-Carlo policy gradient still has high variance

- We can 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

$$\nabla_\theta J(\theta) \approx E_{\pi_\theta}[\nabla_\theta \log \pi_\theta(s, a) Q_w(s, a)]$$

# Actor-Critic

- Monte-Carlo policy gradient still has high variance

- We can 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

$$\nabla_\theta J(\theta) \approx E_{\pi_\theta}[\nabla_\theta \log \pi_\theta(s, a) Q_w(s, a)]$$
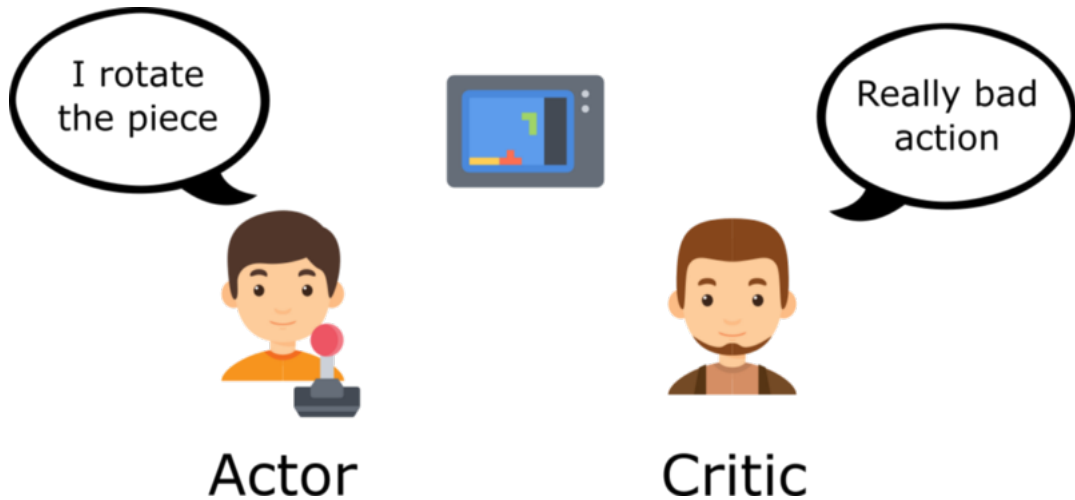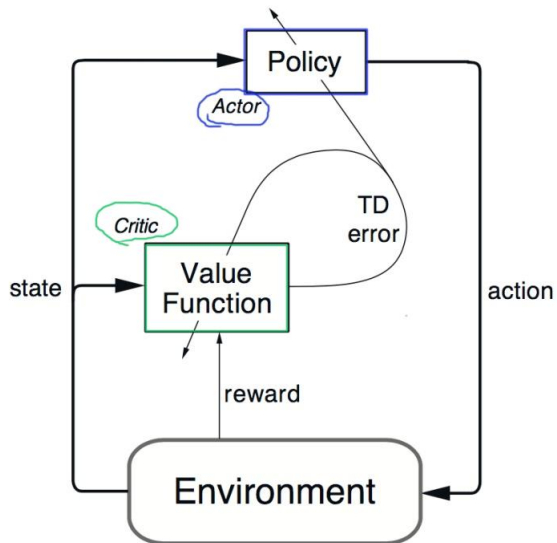$$\Delta\theta = \alpha \nabla_\theta \log \pi_\theta(s, a) Q_w(s, a)$$

- The actor is the policy $\pi_\theta(a|s)$ with parameters $\theta$ which conducts actions in an environment.

- The actor is the policy $\pi_\theta(a|s)$ with parameters $\theta$ which conducts actions in an environment.

- The critic computes value functions to help assist the actor in learning. These are usually the state value, state-action value, or advantage value, denoted as $V(s)$, $Q(s, a)$, and $A(s, a)$, respectively.

# Actor-Critic

- The critic is solving a familiar problem: policy evaluation
- How good is policy $\pi_\theta$ for current parameters $\theta$?

- The critic is solving a familiar problem: policy evaluation

- How good is policy $\pi_\theta$ for current parameters $\theta$?

- To estimate, use any policy evaluation method:

  - Monte-Carlo policy evaluation

  - Temporal-Difference learning

  - Least-squares policy evaluation

- For the true value function $V_{\pi_\theta}(s)$, the TD error $\delta_{\pi_\theta}$

$$\delta_{\pi_\theta} =$$

# Estimating the TD Error

- For the true value function $V_{\pi_\theta}(s)$, the TD error $\delta_{\pi_\theta}$

$$\delta_{\pi_\theta} = r + \gamma V_{\pi_\theta}(s') - V_{\pi_\theta}(s)$$

## Estimating the TD Error

- For the true value function $V_{\pi_\theta}(s)$, the TD error $\delta_{\pi_\theta}$

$$\delta_{\pi_\theta} = r + \gamma V_{\pi_\theta}(s') - V_{\pi_\theta}(s)$$

- is an unbiased estimate of the advantage function

$$\mathbb{E}_{\pi_\theta}[\delta_{\pi_\theta}|s, a] = \mathbb{E}_{\pi_\theta}\left[r + \gamma V_{\pi_\theta}(s')|s, a\right] - V_{\pi_\theta}(s)$$

## Estimating the TD Error

- For the true value function $V_{\pi_\theta}(s)$, the TD error $\delta_{\pi_\theta}$

$$\delta_{\pi_\theta} = r + \gamma V_{\pi_\theta}(s') - V_{\pi_\theta}(s)$$

- is an unbiased estimate of the advantage function

$$\mathbb{E}_{\pi_\theta}[\delta_{\pi_\theta}|s,a] = \mathbb{E}_{\pi_\theta}\left[r + \gamma V_{\pi_\theta}(s')|s,a\right] - V_{\pi_\theta}(s)$$

$$= Q_{\pi_\theta}(s,a) - V_{\pi_\theta}(s)$$

## Estimating the TD Error

- For the true value function $V_{\pi_\theta}(s)$, the TD error $\delta_{\pi_\theta}$

$$\delta_{\pi_\theta} = r + \gamma V_{\pi_\theta}(s') - V_{\pi_\theta}(s)$$

- is an unbiased estimate of the advantage function

$$\begin{aligned}
\mathbb{E}_{\pi_\theta}[\delta_{\pi_\theta}|s, a] &= \mathbb{E}_{\pi_\theta}\left[r + \gamma V_{\pi_\theta}(s')|s, a\right] - V_{\pi_\theta}(s) \\
&= Q_{\pi_\theta}(s, a) - V_{\pi_\theta}(s) \\
&= A_{\pi_\theta}(s, a)
\end{aligned}$$

## Estimating the TD Error

- For the true value function $V_{\pi_\theta}(s)$, the TD error $\delta_{\pi_\theta}$

$$\delta_{\pi_\theta} = r + \gamma V_{\pi_\theta}(s') - V_{\pi_\theta}(s)$$

- is an unbiased estimate of the advantage function

$$\mathbb{E}_{\pi_\theta}[\delta_{\pi_\theta}|s, a] = \mathbb{E}_{\pi_\theta}\left[r + \gamma V_{\pi_\theta}(s')|s, a\right] - V_{\pi_\theta}(s)$$
$$= Q_{\pi_\theta}(s, a) - V_{\pi_\theta}(s)$$
$$= A_{\pi_\theta}(s, a)$$

- So we can use the TD error to compute the policy gradient

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

## Estimating the TD Error

- For the true value function $V_{\pi_\theta}(s)$, the TD error $\delta_{\pi_\theta}$

$$\delta_{\pi_\theta} = r + \gamma V_{\pi_\theta}(s') - V_{\pi_\theta}(s)$$

- is an unbiased estimate of the advantage function

$$\mathbb{E}_{\pi_\theta}[\delta_{\pi_\theta}|s, a] = \mathbb{E}_{\pi_\theta}\left[r + \gamma V_{\pi_\theta}(s')|s, a\right] - V_{\pi_\theta}(s)$$
$$= Q_{\pi_\theta}(s, a) - V_{\pi_\theta}(s)$$
$$= A_{\pi_\theta}(s, a)$$

- So we can use the TD error to compute the policy gradient

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

- In practice we can use an approximate TD error, that requires one set of parameters $w$

$$\delta_w = r + \gamma V_w(s') - V_w(s)$$

# Actor-Critic: Critic (Linear TD(0)) + Actor (policy gradient)

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

Input: a differentiable policy parameterization $\pi(a|s, \boldsymbol{\theta})$
Input: a differentiable state-value function parameterization $\hat{v}(s, \mathbf{w})$
Parameters: step sizes $\alpha^{\boldsymbol{\theta}} > 0$, $\alpha^{\mathbf{w}} > 0$
Initialize policy parameter $\boldsymbol{\theta} \in \mathbb{R}^{d'}$ 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):
        $A \sim \pi(\cdot|S, \boldsymbol{\theta})$
        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 $\hat{v}(S', \mathbf{w}) \doteq 0$)
        $\mathbf{w} \leftarrow \mathbf{w} + \alpha^{\mathbf{w}} \delta \nabla \hat{v}(S, \mathbf{w})$
        $\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} + \alpha^{\boldsymbol{\theta}} I \delta \nabla \ln \pi(A|S, \boldsymbol{\theta})$
        $I \leftarrow \gamma I$
        $S \leftarrow S'$

**REINFORCE with Baseline (episodic), for estimating $\pi_{\boldsymbol{\theta}} \approx \pi_*$**

Input: a differentiable policy parameterization $\pi(a|s, \boldsymbol{\theta})$
Input: a differentiable state-value function parameterization $\hat{v}(s, \mathbf{w})$
Algorithm parameters: step sizes $\alpha^{\boldsymbol{\theta}} > 0$, $\alpha^{\mathbf{w}} > 0$
Initialize policy parameter $\boldsymbol{\theta} \in \mathbb{R}^{d'}$ and state-value weights $\mathbf{w} \in \mathbb{R}^d$ (e.g., to $\mathbf{0}$)

Loop forever (for each episode):
    Generate an episode $S_0, A_0, R_1, \ldots, S_{T-1}, A_{T-1}, R_T$, following $\pi(\cdot|\cdot, \boldsymbol{\theta})$
    Loop for each step of the episode $t = 0, 1, \ldots, T-1$:
        $G \leftarrow \sum_{k=t+1}^{T} \gamma^{k-t-1} R_k$                                    $(G_t)$
        $\delta \leftarrow G - \hat{v}(S_t, \mathbf{w})$
        $\mathbf{w} \leftarrow \mathbf{w} + \alpha^{\mathbf{w}} \delta \nabla \hat{v}(S_t, \mathbf{w})$
        $\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} + \alpha^{\boldsymbol{\theta}} \gamma^t \delta \nabla \ln \pi(A_t|S_t, \boldsymbol{\theta})$

# Advantage Actor Critic (A2C)

- The advantage function can significantly reduce variance of policy gradient

# Advantage Actor Critic (A2C)

- The advantage function can significantly reduce variance of policy gradient
- So the critic should really estimate the advantage function
- For example, by estimating both $V_{\pi_\theta}(s)$ and $Q_{\pi_\theta}(s, a)$

# Advantage Actor Critic (A2C)

- The advantage function can significantly reduce variance of policy gradient

- So the critic should really estimate the advantage function

- For example, by estimating both $V_{\pi_\theta}(s)$ and $Q_{\pi_\theta}(s, a)$

- Using two function approximators and two parameter vectors,

$$V_v(s) \approx V_{\pi_\theta}(s)$$

# Advantage Actor Critic (A2C)

- The advantage function can significantly reduce variance of policy gradient
- So the critic should really estimate the advantage function
- For example, by estimating both $V_{\pi_\theta}(s)$ and $Q_{\pi_\theta}(s, a)$
- Using two function approximators and two parameter vectors,

$$V_v(s) \approx V_{\pi_\theta}(s)$$
$$Q_w(s, a) \approx Q_{\pi_\theta}(s, a)$$

# Advantage Actor Critic (A2C)

- The advantage function can significantly reduce variance of policy gradient

- So the critic should really estimate the advantage function

- For example, by estimating both $V_{\pi_\theta}(s)$ and $Q_{\pi_\theta}(s,a)$

- Using two function approximators and two parameter vectors,

$$V_v(s) \approx V_{\pi_\theta}(s)$$
$$Q_w(s,a) \approx Q_{\pi_\theta}(s,a)$$
$$A(s,a) = Q_w(s,a) - V_v(s)$$

# Advantage Actor Critic (A2C)

- The advantage function can significantly reduce variance of policy gradient

- So the critic should really estimate the advantage function

- For example, by estimating both $V_{\pi_\theta}(s)$ and $Q_{\pi_\theta}(s,a)$

- Using two function approximators and two parameter vectors,

$$V_v(s) \approx V_{\pi_\theta}(s)$$
$$Q_w(s,a) \approx Q_{\pi_\theta}(s,a)$$
$$A(s,a) = Q_w(s,a) - V_v(s)$$

- And updating *both* value functions by e.g. TD learning

# Summary of Policy Gradient Algorithms

- The policy gradient has many equivalent forms

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

- The policy gradient has many equivalent forms

$$\nabla_\theta J(\theta) = \mathbb{E}_{\pi_\theta}[\nabla_\theta \log \pi_\theta(s, a) G_t] \qquad \text{REINFORCE}$$
$$= \mathbb{E}_{\pi_\theta}[\nabla_\theta \log \pi_\theta(s, a) Q_w(s, a)]$$

- The policy gradient has many equivalent forms

$$\nabla_\theta J(\theta) = \mathbb{E}_{\pi_\theta}[\nabla_\theta \log \pi_\theta(s, a) G_t] \qquad \text{REINFORCE}$$
$$= \mathbb{E}_{\pi_\theta}[\nabla_\theta \log \pi_\theta(s, a) Q_w(s, a)] \qquad \text{Q Actor-Critic}$$
$$= \mathbb{E}_{\pi_\theta}[\nabla_\theta \log \pi_\theta(s, a) A_w(s, a)]$$

## Summary of Policy Gradient Algorithms

- The policy gradient has many equivalent forms

$$
\begin{aligned}
\nabla_\theta J(\theta) &= \mathbb{E}_{\pi_\theta}[\nabla_\theta \log \pi_\theta(s, a) G_t] && \text{REINFORCE} \\
&= \mathbb{E}_{\pi_\theta}[\nabla_\theta \log \pi_\theta(s, a) Q_w(s, a)] && \text{Q Actor-Critic} \\
&= \mathbb{E}_{\pi_\theta}[\nabla_\theta \log \pi_\theta(s, a) A_w(s, a)] && \text{Advantage Actor-Critic (A2C)} \\
&= \mathbb{E}_{\pi_\theta}[\nabla_\theta \log \pi_\theta(s, a) \delta]
\end{aligned}
$$

- The policy gradient has many equivalent forms

$$
\begin{aligned}
\nabla_\theta J(\theta) &= \mathbb{E}_{\pi_\theta}[\nabla_\theta \log \pi_\theta(s, a) G_t] && \text{REINFORCE} \\
&= \mathbb{E}_{\pi_\theta}[\nabla_\theta \log \pi_\theta(s, a) Q_w(s, a)] && \text{Q Actor-Critic} \\
&= \mathbb{E}_{\pi_\theta}[\nabla_\theta \log \pi_\theta(s, a) A_w(s, a)] && \text{Advantage Actor-Critic (A2C)} \\
&= \mathbb{E}_{\pi_\theta}[\nabla_\theta \log \pi_\theta(s, a) \delta] && \text{TD Actor-Critic}
\end{aligned}
$$

- Each leads a stochastic gradient ascent algorithm

- Critic uses policy evaluation (e.g. MC or TD learning) to estimate $Q_\pi(s, a)$, $A_\pi(s, a)$ or $V_\pi(s)$.