# Lecture 8:  Policy Gradient

Chonggang Wang

# Outline

- Policy Gradient RL

- Actor-Critic Methods

- Policy Gradient w/ Advantage Function

*materials are modified from David Silver's RL lecture notes

# Outline

- <span style="color:red">Policy Gradient RL</span>

- Actor-Critic Methods

- Policy Gradient w/ Advantage Function

# Lecture 7 – Value Approximation

- Parameterized value function $\hat{v}(s, w)$ and action-value function $\hat{q}(s, a, w)$ using a parameter column vector $w$, to approximate true value functions $V_\pi(s)$ and $Q_\pi(s)$

$$\hat{v}(s, w) \approx V_\pi(s)$$
$$\hat{q}(s, a, w) \approx Q_\pi(s, a)$$

  - How to obtain $\hat{v}(s, w)$ and $\hat{q}(s, a, w)$
    - Stochastic Gradient Descent: $\Delta w = \alpha * (V_\pi(s) - \hat{v}(s, w)) * \nabla_w \hat{v}(s, w)$
    - Calculate the gradient $\nabla_w \hat{v}(s, w)$: Represent a state by a feature vector
      - Linear Approach: $\hat{v}(s, w) = x(s)^T w$
      - Non-Linear Approach: Neural Network (e.g., DQN)

- Use the action-value function to determine/select next action based on a pre-configured policy (e.g., greedy, $\varepsilon$-greedy)

- Question: Can we determine next action directly from a parameterized probability function which only depends on the states: Policy-based RL

$$\pi_\theta(s, a) = p(a|s, \theta)$$

# Why Policy-based RL

- Advantages
  - Action probabilities change more smoothly → better convergence
  - Easy to learn stochastic policies -> better exploration
  - Avoid "argmax( q(s, a) )" → Lower complexity and more effective in high-dimensional or continuous  action spaces

- Disadvantages
  - Value-based approaches could be more efficient for a small number of states and actions
  - Hard to get unbiased estimates of policy gradient through sampling
  - Basic policy gradient approaches (e.g., REINFOCE as a Monte Carlo Method) introduce high variance and lead to a lower convergence speed

# 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 $\pi_\theta$

# Policy Optimization

- Policy based reinforcement learning is an <span style="color:red">optimisation</span> problem
- Find $\theta$ that maximises $J(\theta)$

- Similar to the value based function approximation, we focus on gradient method
  - Gradient is a key to connect neural network with RL algorithms
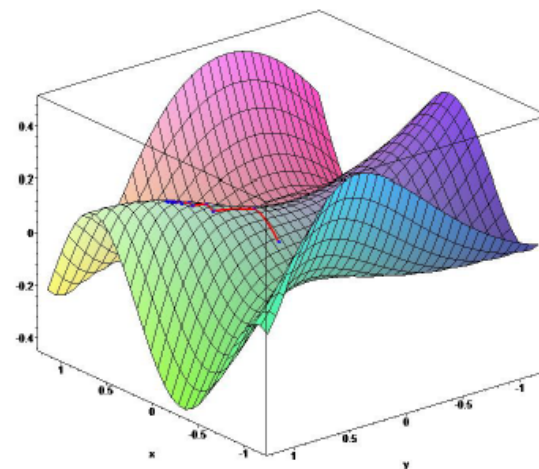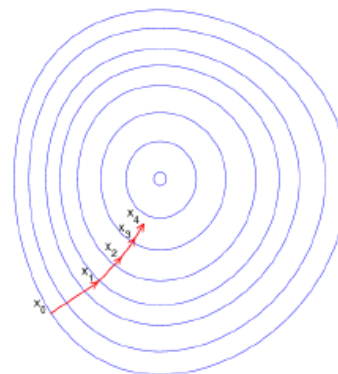- Other approaches are possible

# Gradient Descent (recap)

- Let $J(\mathbf{w})$ be a differentiable function of parameter vector $\mathbf{w}$ , a column vector

- Define the *gradient* of $J(\mathbf{w})$ to be

$$\nabla_{\mathbf{w}} J(\mathbf{w}) = \begin{pmatrix} \frac{\partial J(\mathbf{w})}{\partial \mathbf{w}_1} \\ \vdots \\ \frac{\partial J(\mathbf{w})}{\partial \mathbf{w}_n} \end{pmatrix}$$ Partial Derivatives with respect to $\mathbf{w}$

- To find a local minimum of $J(\mathbf{w})$

- Adjust $\mathbf{w}$ in direction of -ve gradient to reduce $J(w)$ (i.e., the Value Error (VE))

$$\Delta \mathbf{w} = -\frac{1}{2} \alpha \nabla_{\mathbf{w}} J(\mathbf{w})$$

# Gradient Ascent

- $J(\theta)$ is the policy-related performance function to be maximized
- $J(\theta)$ is differentiable with respect to $\theta$, which is a column parameter vector
$$\theta = (\theta_1, \dots, \theta_n)^T$$

- $\nabla_\theta J(\theta)$: The Gradient of $J(\theta)$ can be calculated

$$\nabla_\theta J(\theta) = \begin{pmatrix} \frac{\partial J(\theta)}{\theta_1} \\ . \\ . \\ . \\ \frac{\partial J(\theta)}{\theta_n} \end{pmatrix} \quad \text{Partial Derivatives with respect to } \theta$$

- **Gradient Ascent**
  - To find the optimal parameter vector which maximizes $J(\theta)$

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

# Score Function

- We now compute the policy gradient *analytically*
- Assume policy $\pi_\theta$ is differentiable whenever it is non-zero and we know the gradient $\nabla_\theta \pi_\theta(s, a)$
- Likelihood ratios exploit the following identity

$$\nabla_\theta \pi_\theta(s, a) = \pi_\theta(s, a) \frac{\nabla_\theta \pi_\theta(s, a)}{\pi_\theta(s, a)}$$
$$= \pi_\theta(s, a) \nabla_\theta \ln \pi_\theta(s, a)$$

- The score function is $\nabla_\theta \ln \pi_\theta(s, a)$

# Example: Softmax Policy

- We will use a softmax policy as a running example
- Weight actions using linear combination of features $\phi(s,a)^\top \theta$
- Probability of action is proportional to exponentiated weight

$$\pi_\theta(s,a) = p(a|s,\theta) = \frac{e^{\phi(s,a)^T \theta}}{\sum_a e^{\phi(s,a)^T \theta}}$$

- The score function is

$$\nabla_\theta \ln \pi_\theta(s,a) = \phi(s,a) - \mathbb{E}_{\pi_\theta}[\phi(s,\cdot)]$$

# Example: Softmax Policy

- Proof:

$$\nabla \ln \pi_\theta(s,a) = \nabla \ln \frac{e^{\phi(s,a)^T \theta}}{\sum_a e^{\phi(s,a)^T \theta}} \qquad \text{--Eq. (1)}$$

$$\nabla \ln\left(\frac{x}{y}\right) = \nabla \ln(x) - \nabla\ln(y)$$

$$= \nabla \phi(s,a)^T \theta - \nabla \ln\left(\sum_a e^{\phi(s,a)^T \theta}\right) \qquad \text{--Eq. (2)}$$

$$\nabla \ln(g(x)) = \frac{\nabla g(x)}{g(x)}$$

$$= \phi(s,a) - \frac{\sum_a (e^{\phi(s,a)^T \theta} \phi(s,a))}{\sum_a e^{\phi(s,a)^T \theta}} \qquad \text{--Eq. (3)}$$

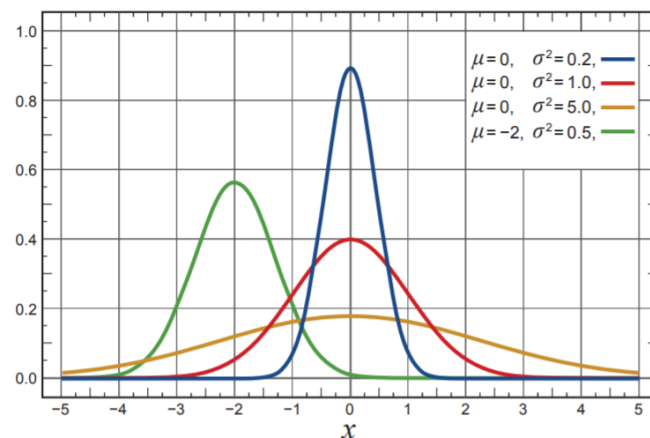$$= \phi(s,a) - \sum_a \frac{e^{\phi(s,a)^T \theta}}{\sum_a e^{\phi(s,a)^T \theta}} \phi(s,a) \qquad \text{--Eq. (4)}$$

$$= \phi(s,a) - \sum_a \pi_\theta(s,a) \phi(s,a) \qquad \text{--Eq. (5)}$$

$$= \phi(s,a) - \mathbb{E}_{\pi_\theta}[\phi(s,\cdot)] \qquad \text{--Eq. (6)}$$

# Example: Gaussian Policy

- In continuous action spaces, a Gaussian policy is natural
- Mean is a linear combination of state features $\mu(s) = \phi(s)^\top \theta$
- Variance may be fixed $\sigma^2$, or can also parametrised
- Policy is Gaussian, $a \sim \mathcal{N}(\mu(s), \sigma^2)$   $\pi_\theta(a|s, \theta) = \dfrac{1}{\sigma\sqrt{2\pi}} e^{\left(-\frac{(a - \mu(s, \theta))^2}{2\sigma^2}\right)}$
- The score function is

$$\nabla_\theta \ln \pi_\theta(s, a) = \frac{(a - \mu(s))\phi(s)}{\sigma^2}$$



Legend: $\mu=0$, $\sigma^2=0.2$; $\mu=0$, $\sigma^2=1.0$; $\mu=0$, $\sigma^2=5.0$; $\mu=-2$, $\sigma^2=0.5$

# One-Step MDP

- 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

$$J(\theta) = \mathbb{E}_{\pi_\theta}[r]$$

$$= \sum_{s \in \mathcal{S}} d(s) \sum_{a \in \mathcal{A}} \pi_\theta(s,a) \mathcal{R}_{s,a}$$

$$\nabla_\theta J(\theta) = \sum_{s \in \mathcal{S}} d(s) \sum_{a \in \mathcal{A}} \pi_\theta(s,a) \nabla_\theta \ln \pi_\theta(s,a) \mathcal{R}_{s,a}$$

$$= \mathbb{E}_{\pi_\theta}[\nabla_\theta \ln \pi_\theta(s,a) r]$$

# Generalized MDP

- The policy gradient theorem generalises 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

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

$$\nabla_\theta J(\theta) = \mathbb{E}_{\pi_\theta} [\nabla_\theta \ln \pi_\theta(s, a) \ Q^{\pi_\theta}(s, a)]$$

Policy Gradient Theorem

# Monte-Carlo Policy Gradient (REINFORCE)

- Update parameters by stochastic gradient ascent
- Using policy gradient theorem
- Using <u>return $G_t$</u> as an unbiased sample of $Q^{\pi_\theta}(s_t, a_t)$

High Variance

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

Policy Gradient

```
1   function REINFORCE
2       Initialise θ arbitrarily
3       for each episode {s_1, a_1, r_2, ..., s_{T-1}, a_{T-1}, r_T} ~ π_θ do
4           for t = 1 to T − 1 do
5               θ ← θ + α∇_θ ln π_θ(s_t, a_t) G_t
6           end for
7       end for
8       return θ
9   end function
```

# Outline

- Policy Gradient RL

- <span style="color:red">Actor-Critic Methods</span>

- Policy Gradient w/ Advantage Function

*materials are modified from David Silver's RL lecture notes

# 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

$$\nabla_\theta J(\theta) \approx \mathbb{E}_{\pi_\theta} [\nabla_\theta \ln \pi_\theta(s, a) \; Q_w(s, a)]$$
$$\Delta\theta = \alpha \nabla_\theta \ln \pi_\theta(s, a) \; Q_w(s, a)$$

# Example: TD(0)-based Actor-Critic Algorithm

- Simple actor-critic algorithm based on action-value critic
- Using linear value fn approx. $Q_w(s, a) = \phi(s, a)^\top w$

  Critic Updates $w$ by linear TD(0)

  Actor Updates $\theta$ by policy gradient

1— **function** $\mathrm{QAC}$
2—     Initialise $s$, $\theta$
3—     Sample $a \sim \pi_\theta$
4—     **for** each step **do**
5—         Sample reward $r = \mathcal{R}_s^a$; sample transition $s' \sim \mathcal{P}_{s,\cdot}^a$
6—         Sample action $a' \sim \pi_\theta(s', a')$
7—         $\delta = r + \gamma Q_w(s', a') - Q_w(s, a)$
8—         $\theta = \theta + \alpha \nabla_\theta \ln \pi_\theta(s, a) Q_w(s, a)$
9—         $w \leftarrow w + \beta \delta \phi(s, a)$ ⟵ $w = w + \beta * \delta * (\nabla_w Q_w(s, a))$
10—        $a \leftarrow a', s \leftarrow s'$
11—    **end for**
12— **end function**

# Problem: Bias in Actor-Critic Algorithms

- Approximating the policy gradient introduces bias
- A biased policy gradient may not find the right solution


- Luckily, if we choose value function approximation carefully
- Then we can avoid introducing any bias

  i.e. We can still follow the *exact* policy gradient

# Compatible Function Approximation Theorem*

If the following two conditions are satisfied:

Value function approximator is *compatible* to the policy

$$\nabla_w Q_w(s, a) = \nabla_\theta \ln \pi_\theta(s, a)$$

Value function parameters $w$ minimise the mean-squared error

$$\varepsilon = \mathbb{E}_{\pi_\theta} \left[ (Q^{\pi_\theta}(s, a) - Q_w(s, a))^2 \right]$$

$$\frac{\partial \varepsilon}{\partial w} = 0$$

Then the policy gradient is exact,

$$\nabla_\theta J(\theta) = \mathbb{E}_{\pi_\theta} \left[ \nabla_\theta \ln \pi_\theta(s, a) \, Q_w(s, a) \right]$$

*R. Sutton, et al. "Policy gradient methods for reinforcement learning with function approximation", 2000

# Outline

- Policy Gradient RL

- Actor-Critic Methods

- Policy Gradient w/ Advantage Function

# Reducing Variance Using a Baseline

$$\nabla_\theta J(\theta) = \mathbb{E}_{\pi_\theta} \left[ \nabla_\theta \ln \pi_\theta(s, a) \, Q^{\pi_\theta}(s, a) \right]$$ (Without a Baseline)

- We subtract a baseline function $B(s)$ from the policy gradient
- This can reduce variance, without changing expectation

$$\mathbb{E}_{\pi_\theta} \left[ \nabla_\theta \ln \pi_\theta(s, a) B(s) \right] = \sum_{s \in \mathcal{S}} d^{\pi_\theta}(s) \sum_a \nabla_\theta \pi_\theta(s, a) B(s)$$

*B(s)* is action-independent

$$= \sum_{s \in \mathcal{S}} d^{\pi_\theta} B(s) \nabla_\theta \sum_{a \in \mathcal{A}} \pi_\theta(s, a)$$

$$= 0$$

- A good baseline is the state value function $B(s) = V^{\pi_\theta}(s)$
- So we can rewrite the policy gradient using the advantage function $A^{\pi_\theta}(s, a)$

$$A^{\pi_\theta}(s, a) = Q^{\pi_\theta}(s, a) - V^{\pi_\theta}(s)$$

$$\nabla_\theta J(\theta) = \mathbb{E}_{\pi_\theta} \left[ \nabla_\theta \ln \pi_\theta(s, a) \, A^{\pi_\theta}(s, a) \right]$$ (Using a Baseline)

"Variance Reduction for Policy Gradient with Action-Dependent Factorized Baselines " - ICLR 2018

# Estimating the Advantage Function

- 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

# Estimating the Advantage Function (cont.)

- 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

$$\mathbb{E}_{\pi_\theta}\left[\delta^{\pi_\theta}|s, a\right] = \mathbb{E}_{\pi_\theta}\left[r + \gamma V^{\pi_\theta}(s')|s, a\right] - V^{\pi_\theta}(s)$$
$$= Q^{\pi_\theta}(s, a) - V^{\pi_\theta}(s)$$
$$= A^{\pi_\theta}(s, a)$$

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

$$\nabla_\theta J(\theta) = \mathbb{E}_{\pi_\theta}\left[\nabla_\theta \boldsymbol{ln}\, \pi_\theta(s, a)\, \delta^{\pi_\theta}\right]$$

- In practice we can use an approximate TD error

$$\delta_v = r + \gamma V_v(s') - V_v(s)$$

- This approach only requires one set of critic parameters $v$

# Critic at Different Time-Scales

- Critic can estimate value function $V_v(s)$ from many targets at different time-scales

  For MC, the target is the return $v_t$

  $$\Delta V = \alpha(\, G_t - V_v(s)\,)\phi(s)$$

  linear approximation

  For TD(0), the target is the TD target $r + \gamma V(s')$

  $$\Delta V = \alpha(r + \gamma V(s') - V_v(s))\phi(s)$$

  For forward-view TD($\lambda$), the target is the $\lambda$-return $v_t^\lambda$

  $$\Delta V = \alpha(G_t^\lambda - V_v(s))\phi(s)$$

  For backward-view TD($\lambda$), we use eligibility traces

  $$\delta_t = r_{t+1} + \gamma V(s_{t+1}) - V(s_t)$$
  $$e_t = \gamma \lambda e_{t-1} + \phi(s_t)$$
  $$\Delta V = \alpha \delta_t e_t$$

# Actor at Different Time-Scales

- The policy gradient can also be estimated at many time-scales

$$\nabla_\theta J(\theta) = \mathbb{E}_{\pi_\theta} [\nabla_\theta \, ln \, \pi_\theta(s, a) \, A^{\pi_\theta}(s, a)]$$

- Monte-Carlo policy gradient uses error from complete return

$$\Delta\theta = \alpha( G_t - V_v(s_t))\nabla_\theta \, ln \, \pi_\theta(s_t, a_t)$$

- Actor-critic policy gradient uses the one-step TD error

$$\Delta\theta = \alpha(r + \gamma V_v(s_{t+1}) - V_v(s_t))\nabla_\theta \, ln \, \pi_\theta(s_t, a_t)$$

# Policy Gradient with Eligibility Traces

- Just like forward-view TD($\lambda$), we can mix over time-scales

$$\Delta\theta = \alpha(G_t^\lambda - V_v(s_t))\nabla_\theta \ ln \ \pi_\theta(s_t, a_t)$$

where $G_t^\lambda - V_v(s_t)$ is a biased estimate of advantage fn
- Like backward-view TD($\lambda$), we can also use eligibility traces

By equivalence with TD($\lambda$), substituting $\phi(s) = \nabla_\theta \ ln \ \pi_\theta(s, a)$

$$\delta = r_{t+1} + \gamma V_v(s_{t+1}) - V_v(s_t)$$
$$e_{t+1} = \lambda e_t + \nabla_\theta \ ln \ \pi_\theta(s, a)$$
$$\Delta\theta = \alpha\delta e_t$$

- This update can be applied online, to incomplete sequences

# Summary of Policy Gradient Algorithms

- The policy gradient has many equivalent forms

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

$$= \mathbb{E}_{\pi_\theta}\left[\nabla_\theta \ln \pi_\theta(s,a)\, Q^w(s,a)\right] \qquad \text{Q Actor-Critic}$$

$$= \mathbb{E}_{\pi_\theta}\left[\nabla_\theta \ln \pi_\theta(s,a)\, A^w(s,a)\right] \qquad \text{Advantage Actor-Critic}$$

$$= \mathbb{E}_{\pi_\theta}\left[\nabla_\theta \ln \pi_\theta(s,a)\, \delta\right] \qquad \text{TD Actor-Critic}$$

$$= \mathbb{E}_{\pi_\theta}\left[\nabla_\theta \ln \pi_\theta(s,a)\, \delta e\right] \qquad \text{TD}(\lambda)\text{ Actor-Critic}$$

- Each leads to a stochastic gradient ascent algorithm.