

# Lecture 8:

# Policy Gradient Algorithms

**B. Ravindran**

# Solution Methods

- ❑ Value-based Methods
  - ❑ Q-learning
  - ❑ SARSA
  - ❑ TD( $\lambda$ )
  - ❑ Actor-Critic
- ❑ Policy Search
  - ❑ Policy Gradient Methods
  - ❑ Evolutionary algorithms
- ❑ Model based methods
  - ❑ Stochastic Dynamic Programming
  - ❑ Bayesian approaches

# Solution Methods

- ❑ Value-based Methods
  - ❑ Q-learning
  - ❑ SARSA
  - ❑ TD( $\lambda$ )
  - ❑ Actor-Critic
- ❑ Policy Search
  - ❑ Policy Gradient Methods
  - ~~❑ Evolutionary algorithms~~
- ❑ Model based methods
  - ❑ Stochastic Dynamic Programming
  - ~~❑ Bayesian approaches~~

# Policy Search Methods

- ❑ Policy search: Instead of maintaining estimates of value functions, search in the space of policies
- ❑ **Why?**
  - ❑ Simpler description
  - ❑ Better convergence
  - ❑ Robust to partial observability
  - ❑ Continuous action space
- ❑ Direct policy search - Genetic algorithms
- ❑ Policy Gradient Approaches

# Policy Gradient Methods

- ❑ Policy depends on some parameters  $\theta$ 
  - ❑ Action preferences
  - ❑ Mean and variance
  - ❑ Weights of a neural network
- ❑ **Idea:** Modify policy parameters directly instead of estimating the action values

- ❑ Maximize:  $J(\theta) = \mathbb{E}[G_t]$   $J(\theta) = \mathbb{E}[r_t]$

Simplified Setting

Immediate Reward or  
Multi-arm bandits

- ❑  $\theta$  update:  $\theta \leftarrow \theta + \alpha \nabla J(\theta)$

# Stochastic Gradient Ascent

- We compute the gradient of the performance  $J(\theta)$  w.r.t the parameters  $\theta$

$$J(\theta) = E(r_t) = \sum_a q_*(a) \pi_\theta(a)$$

$$\nabla J(\theta) = \sum_a q_*(a) \nabla \pi_\theta(a)$$

$$= \sum_a q_*(a) \frac{\nabla \pi_\theta(a)}{\pi_\theta(a)} \pi_\theta(a) = \mathbb{E}_{\pi(\theta)} [q_*(a) \nabla \ln \pi_\theta(a)]$$



Bandit Setting  
(Immediate Reward)

- Estimate the gradient from  $N$  samples

$$\hat{\nabla} J(\theta) = \frac{1}{N} \sum_{i=1}^N r_i \underbrace{\frac{\nabla \pi_\theta(a_i)}{\pi_\theta(a_i)}}_{\text{Likelihood Ratio}}$$

# REINFORCE

- REward Increment = Non-negative Factor  $\times$  Offset Reinforcement  $\times$  Characteristic Eligibility (*Williams '92*)

- Incremental version

$$\Delta \theta_t = \alpha r_t \frac{\nabla \pi_{\theta}(a_t)}{\pi_{\theta}(a_t)}$$

$$\Delta \theta_t = \alpha r_t \frac{\partial \ln \pi_{\theta}(a_t)}{\partial \theta}$$



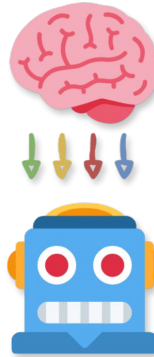
Bandit Setting  
(Immediate Reward)

- REINFORCE with baseline

$$\Delta \theta_t = \alpha (r_t - b_t) \frac{\partial \ln \pi_{\theta}(a_t)}{\partial \theta}$$

Reinforcement  
Baseline

Characteristic  
Eligibility



# Policy Gradient Theorem



# Policy Gradient Theorem

- In the episodic case, we define the performance by assuming that every episode starts from state  $s_0$  (non-random), as follows:

$$J(\theta) \doteq v_{\pi_\theta}(s_0)$$

where  $v_{\pi_\theta}(s_0)$  is the true value function given a parameterized policy  $\pi_\theta$

- **The Policy Gradient Theorem:** The gradient of the performance can be expressed in terms of the gradient of the policy, as follows

$$\nabla J(\theta) \propto \sum_s \mu(s) \sum_a q_\pi(s, a) \nabla \pi(a|s, \theta)$$

# Policy Gradient Theorem

- ❑ Provides an analytic expression for the gradient of performance with respect to the policy parameter  $\theta$
- ❑ We begin the proof by expressing the gradient of the state-value function in terms of the action-value function

$$\begin{aligned}
 \nabla_{\theta} v_{\pi}(s) &= \nabla_{\theta} \left( \sum_{a \in \mathcal{A}} \pi_{\theta}(a | s) q_{\pi}(s, a) \right) \quad \text{Derivative product rule} \\
 &= \sum_{a \in \mathcal{A}} (\nabla_{\theta} \pi_{\theta}(a | s) q_{\pi}(s, a) + \pi_{\theta}(a | s) \nabla_{\theta} q_{\pi}(s, a)) \quad \text{Write } Q_{\pi} \text{ as } V_{\pi} \\
 &= \sum_{a \in \mathcal{A}} \left( \nabla_{\theta} \pi_{\theta}(a | s) q_{\pi}(s, a) + \pi_{\theta}(a | s) \nabla_{\theta} \sum_{s', r} p(s', r | s, a) (r + v_{\pi}(s')) \right) \quad P \text{ is not a function of } \theta \\
 &= \sum_{a \in \mathcal{A}} \left( \nabla_{\theta} \pi_{\theta}(a | s) q_{\pi}(s, a) + \pi_{\theta}(a | s) \sum_{s', r} p(s', r | s, a) \nabla_{\theta} v_{\pi}(s') \right) \\
 &= \sum_{a \in \mathcal{A}} \left( \nabla_{\theta} \pi_{\theta}(a | s) q_{\pi}(s, a) + \pi_{\theta}(a | s) \sum_{s'} p(s' | s, a) \nabla_{\theta} v_{\pi}(s') \right) \quad \sum_r P(r|s, a) = 1
 \end{aligned}$$

# Policy Gradient Theorem

$$\nabla_{\theta} v_{\pi}(s) = \sum_{a \in \mathcal{A}} \left( \nabla_{\theta} \pi_{\theta}(a | s) q_{\pi}(s, a) + \pi_{\theta}(a | s) \sum_{s'} p(s' | s, a) \nabla_{\theta} v_{\pi}(s') \right)$$

$$\nabla_{\theta} v_{\pi}(s) = \sum_{a \in \mathcal{A}} \left( \nabla_{\theta} \pi_{\theta}(a | s) q_{\pi}(s, a) + \pi_{\theta}(a | s) \sum_{s'} p(s' | s, a) \left( \sum_{a' \in \mathcal{A}} \left( \nabla_{\theta} \pi_{\theta}(a' | s') q_{\pi}(s', a') + \pi_{\theta}(a' | s') \sum_{s''} p(s'' | s', a') \nabla_{\theta} v_{\pi}(s'') \right) \right) \right)$$

-----

Unroll  $\nabla_{\theta} V_{\pi}(s')$

# Policy Gradient Theorem

$$s \xrightarrow{a \sim \pi_\theta(\cdot|s)} s' \xrightarrow{a \sim \pi_\theta(\cdot|s')} s'' \xrightarrow{a \sim \pi_\theta(\cdot|s'')} \dots$$

Probability of transitioning from  $s$  to  $s$  in 0 steps while following  $\pi_\theta$

$$\rho_\pi(s \rightarrow s, k = 0) = 1$$

Probability of transitioning from  $s$  to  $s'$  in 1 step while following  $\pi_\theta$

$$\rho_\pi(s \rightarrow s', k = 1) = \sum_a \pi_\theta(a|s) p(s'|s, a)$$

Probability of transitioning from  $s$  to  $s'$  in  $k$  steps and from  $s'$  to  $x$  in 1 step while following  $\pi_\theta$

$$\rho_\pi(s \rightarrow x, k + 1) = \sum_{s'} \rho_\pi(s \rightarrow s', k) \rho_\pi(s' \rightarrow x, 1)$$

# Policy Gradient Theorem

$$\nabla_{\theta} v_{\pi}(s) = \sum_{a \in \mathcal{A}} \left( \nabla_{\theta} \pi_{\theta}(a | s) q_{\pi}(s, a) + \pi_{\theta}(a | s) \sum_{s'} p(s' | s, a) \nabla_{\theta} v_{\pi}(s') \right)$$

$$\nabla_{\theta} v_{\pi}(s) = \sum_{a \in \mathcal{A}} \left( \nabla_{\theta} \pi_{\theta}(a | s) q_{\pi}(s, a) + \pi_{\theta}(a | s) \sum_{s'} p(s' | s, a) \right) \left( \nabla_{\theta} v_{\pi}(s') \right)$$

← Unroll  $\nabla_{\theta} V_{\pi}(s')$

$$\sum_{a' \in \mathcal{A}} \left( \nabla_{\theta} \pi_{\theta}(a' | s') q_{\pi}(s', a') + \pi_{\theta}(a' | s') \sum_{s''} p(s'' | s', a') \nabla_{\theta} v_{\pi}(s'') \right)$$

-----

$$\phi(s) = \sum_{a \in \mathcal{A}} \nabla_{\theta} \pi_{\theta}(a | s) q_{\pi}(s, a)$$

Distribute  
outer-most  $\sum_a$

$$\sum_{s'} \rho_{\pi}(s \rightarrow s', k=1)$$

$$\nabla_{\theta} v_{\pi}(s) = \left( \sum_{a \in \mathcal{A}} \nabla_{\theta} \pi_{\theta}(a | s) q_{\pi}(s, a) + \sum_{s'} \sum_{a \in \mathcal{A}} \pi_{\theta}(a | s) p(s' | s, a) \right) \left( \sum_{a' \in \mathcal{A}} \left( \nabla_{\theta} \pi_{\theta}(a' | s') q_{\pi}(s', a') + \pi_{\theta}(a' | s') \sum_{s''} p(s'' | s', a') \nabla_{\theta} v_{\pi}(s'') \right) \right)$$

$\phi(s)$                        $\phi(s')$

# Policy Gradient Theorem

$$\nabla_{\theta} v_{\pi}(s) = \left( \sum_{a \in \mathcal{A}} \nabla_{\theta} \pi_{\theta}(a | s) q_{\pi}(s, a) + \sum_{s'} \sum_{a \in \mathcal{A}} \pi_{\theta}(a | s) p(s' | s, a) \right. \\ \left. \sum_{a' \in \mathcal{A}} \left( \nabla_{\theta} \pi_{\theta}(a' | s') q_{\pi}(s', a') + \pi_{\theta}(a' | s') \sum_{s''} p(s'' | s', a') \nabla_{\theta} v_{\pi}(s'') \right) \right) \rho_{\pi}(s \rightarrow s', k=1) \phi(s')$$

$$\nabla_{\theta} v_{\pi}(s) = \phi(s) + \sum_{s'} \rho_{\pi}(s \rightarrow s', 1) [\phi(s') + \sum_{s''} \rho_{\pi}(s' \rightarrow s'', 1) \nabla_{\theta} v_{\pi}(s'')]$$

$$\nabla_{\theta} v_{\pi}(s) = \phi(s) + \sum_{s'} \rho_{\pi}(s \rightarrow s', 1) \phi(s') + \sum_{s''} \rho_{\pi}(s' \rightarrow s'', 2) \nabla_{\theta} v_{\pi}(s'')$$

Distribute  $\Sigma_{s'}$

$$\nabla_{\theta} v_{\pi}(s) = \sum_{x \in \mathcal{S}} \sum_{k=0}^{\infty} \rho_{\pi}(s \rightarrow x, k) \phi(x)$$

Repeatedly unroll till  $\infty$

# Policy Gradient Theorem

$$\nabla J(\theta) = \nabla_{\theta} v_{\pi}(s_0) \quad (\text{Remember})$$

$$\nabla_{\theta} v_{\pi}(s_0) = \sum_{s \in \mathcal{S}} \left( \sum_{k=0}^{\infty} \rho_{\pi}(s_0 \rightarrow s, k) \right) \sum_{a \in \mathcal{A}} \nabla_{\theta} \pi_{\theta}(a|s) q_{\pi}(s, a)$$

$$\nabla_{\theta} v_{\pi}(s_0) = \sum_{s \in \mathcal{S}} \eta(s) \sum_{a \in \mathcal{A}} \nabla_{\theta} \pi_{\theta}(a|s) q_{\pi}(s, a)$$

Introduce the probability version of  $\eta(s)$

$$\nabla_{\theta} v_{\pi}(s_0) = \sum_{s' \in \mathcal{S}} \eta(s') \sum_{s \in \mathcal{S}} \frac{\eta(s)}{\sum_{x \in \mathcal{S}} \eta(x)} \sum_{a \in \mathcal{A}} \nabla_{\theta} \pi_{\theta}(a|s) q_{\pi}(s, a)$$

$$\nabla_{\theta} v_{\pi}(s_0) = \sum_{s' \in \mathcal{S}} \eta(s') \sum_{s \in \mathcal{S}} \mu(s) \sum_{a \in \mathcal{A}} \nabla_{\theta} \pi_{\theta}(a|s) q_{\pi}(s, a)$$

$$\nabla_{\theta} v_{\pi}(s_0) \propto \sum_{s \in \mathcal{S}} \mu(s) \sum_{a \in \mathcal{A}} \nabla_{\theta} \pi_{\theta}(a|s) q_{\pi}(s, a)$$

$\eta(s')$  is a constant,  
 = 1 (in continuing task)  
 = avg episodic length (o.w)

# REINFORCE: MC Policy Gradient

$$\begin{aligned}\nabla J(\boldsymbol{\theta}) &\propto \sum_s \mu(s) \sum_a q_\pi(s, a) \nabla \pi(a|s, \boldsymbol{\theta}) \\&= \mathbb{E}_\pi \left[ \sum_a q_\pi(S_t, a) \nabla \pi(a|S_t, \boldsymbol{\theta}) \right] && \text{Expectation over the visited states while following } \Pi \\&= \mathbb{E}_\pi \left[ q_\pi(S_t, A_t) \frac{\nabla \pi(A_t|S_t, \boldsymbol{\theta})}{\pi(A_t|S_t, \boldsymbol{\theta})} \right] && \begin{array}{l} \text{Multiply and divide by } \Pi(a|s, \boldsymbol{\theta}) + \\ \text{Expectation over the actions taken while following } \Pi \end{array} \\&= \mathbb{E}_\pi \left[ G_t \frac{\nabla \pi(A_t|S_t, \boldsymbol{\theta})}{\pi(A_t|S_t, \boldsymbol{\theta})} \right]\end{aligned}$$

The gradient update is

- proportional to the return
- inversely proportional to the action probability

$$\boldsymbol{\theta}_{t+1} \doteq \boldsymbol{\theta}_t + \alpha G_t \frac{\nabla \pi(A_t|S_t, \boldsymbol{\theta}_t)}{\pi(A_t|S_t, \boldsymbol{\theta}_t)}$$



# REINFORCE: MC PG Control

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  $t = 0, 1, \dots, T - 1$ :

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

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

- ❑ The algorithm computes an unbiased estimate of the gradient
- ❑ Can be very slow due to high variance in the estimates
- ❑ Variance is related to the “recurrence time” or the episode length
- ❑ In large state spaces, the variance becomes unacceptably high

# REINFORCE w/ Baseline

- ❑ The policy gradient theorem can be generalized to include a comparison of the action value to an arbitrary baseline  $b(s)$

$$\nabla J(\boldsymbol{\theta}) \propto \sum_s \mu(s) \sum_a \left( q_\pi(s, a) - b(s) \right) \nabla \pi(a|s, \boldsymbol{\theta})$$

- ❑ Baseline - should be a function that is independent on the action  $a$

$$\sum_a b(s) \nabla \pi(a|s, \boldsymbol{\theta}) = b(s) \nabla \sum_a \pi(a|s, \boldsymbol{\theta}) = b(s) \nabla 1 = 0$$

- ❑ Update rule of REINFORCE with baseline

$$\boldsymbol{\theta}_{t+1} \doteq \boldsymbol{\theta}_t + \alpha \left( G_t - b(S_t) \right) \frac{\nabla \pi(A_t|S_t, \boldsymbol{\theta}_t)}{\pi(A_t|S_t, \boldsymbol{\theta}_t)}$$

(Doesn't change the expected value, but has an effect on the variance)

# Actor-Critic Methods

- ❑ Actor-Critic methods learn both a policy and a state-value function simultaneously
- ❑ The policy is referred to as the actor that suggests actions given a state
- ❑ The value function is referred to as the critic. It evaluates actions taken by the actor based on the given policy

$$\begin{aligned}\boldsymbol{\theta}_{t+1} &\doteq \boldsymbol{\theta}_t + \alpha (G_t - \hat{v}(S_t, \mathbf{w})) \frac{\nabla \pi(A_t | S_t, \boldsymbol{\theta}_t)}{\pi(A_t | S_t, \boldsymbol{\theta}_t)} \\ &= \boldsymbol{\theta}_t + \alpha (R_{t+1} + \gamma \hat{v}(S_{t+1}, \mathbf{w}) - \hat{v}(S_t, \mathbf{w})) \frac{\nabla \pi(A_t | S_t, \boldsymbol{\theta}_t)}{\pi(A_t | S_t, \boldsymbol{\theta}_t)} \\ &= \boldsymbol{\theta}_t + \alpha \delta_t \frac{\nabla \pi(A_t | S_t, \boldsymbol{\theta}_t)}{\pi(A_t | S_t, \boldsymbol{\theta}_t)}\end{aligned}$$

# One-step Actor-Critic Algorithm

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  $\hat{v}(S', \mathbf{w}) \doteq 0$ )  
         $\mathbf{w} \leftarrow \mathbf{w} + \alpha^\mathbf{w} \delta \nabla \hat{v}(S, \mathbf{w})$   
         $\theta \leftarrow \theta + \alpha^\theta I \delta \nabla \ln \pi(A|S, \theta)$   
         $I \leftarrow \gamma I$   
         $S \leftarrow S'$

(This is a fully online, incremental algorithm, with states, actions, and rewards processed as they occur and then never revisited again)