

# Actor Critic Methods: From Paper to Code

Review of Fundamental Concepts

# Agent, Environment, Action



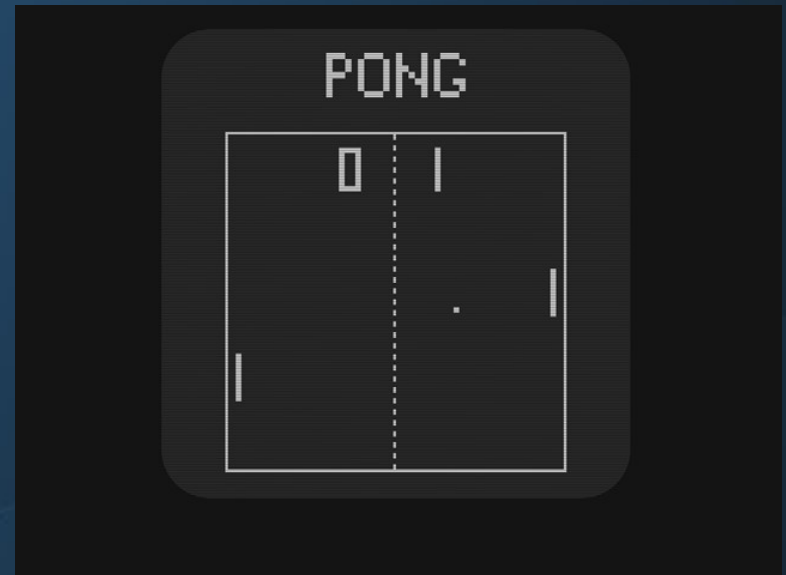
# Markov Decision Process

State depends only on  
previous state and action

Markov Decision Process



$(S_1, A_1, R_1, S_2, A_2, R_2, \dots)$



# Episodic Returns



These states have value



Present rewards worth more

$$G_t = R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} = \sum_{k=0}^{\infty} \gamma^k R_{t+k+1}$$

Sum of discounted rewards → Episode return

# Reward Discounting

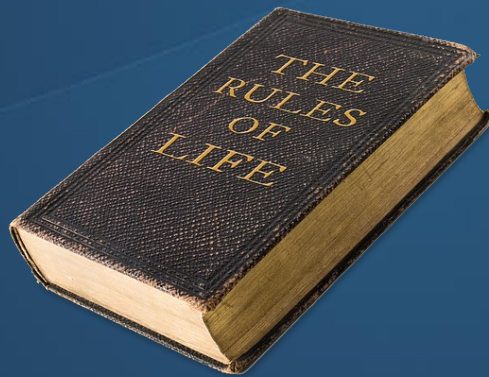
$$G_t = R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \gamma^3 R_{t+4} + \dots$$

$$G_t = R_{t+1} + \gamma (R_{t+2} + \gamma R_{t+3} + \gamma^2 R_{t+4} + \dots)$$

$$G_t = R_{t+1} + \gamma G_{t+1}$$

But wait... how can we know future rewards?

# The Agent's Policy



Mapping of states to actions



Can be probabilistic

# $\Pi$

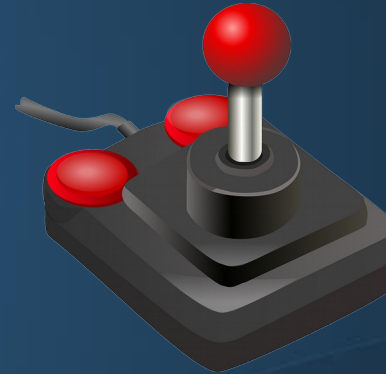
# Value and Action Value Functions

$$v_{\pi}(s) = E_{\pi}[G_t | S_t = s] = E_{\pi}\left[\sum_{k=0}^{\infty} \gamma^k R_{t+k+1} | S_t = s\right] \text{ for all } s \in S$$

$$q_{\pi}(s, a) = E_{\pi}[G_t | S_t = s, A_t = a] = E_{\pi}\left[\sum_{k=0}^{\infty} \gamma^k R_{t+k+1} | S_t = s, A_t = a\right] \text{ for all } s \in S$$

# Learning from Experience

Interact with environment



Keep track of rewards

Monte Carlo Methods



# The Bellman Equation

$$v_{\pi}(s) = E_{\pi}[G_t | S_t = s] = E_{\pi}\left[\sum_{k=0}^{\infty} \gamma^k R_{t+k+1} | S_t = s\right] \text{ for all } s \in S$$

$$G_t = R_{t+1} + \gamma G_{t+1}$$

$$v_{\pi}(s) = E_{\pi}[R_{t+1} + \gamma G_{t+1} | S_t = s]$$

$$v_{\pi}(s) = \sum_a \pi(a, s) \sum_{s'} \sum_r p(s', r | s, a) [r + \gamma E_{\pi}[G_{t+1} | S_{t+1} = s']]$$

$$v_{\pi}(s) = \sum_a \pi(a, s) \sum_{s', r} p(s', r | s, a) [r + \gamma v_{\pi}(s')]$$

## Bellman Equation

$$q_{\pi}(s, a) = \sum_{s', r} p(s', r | s, a) [r + \gamma \sum_{a'} \pi(a' | s') q_{\pi}(s', a')]$$

# Optimal Policies



Compare policies

$$F: I \rightarrow \mathbb{R}, x \mapsto \int_a^x f(t) dt$$
$$\int_a^b f(x) dx = F(b) - F(a)$$

Known dynamics  $\rightarrow$  Model based



Unknown dynamics  $\rightarrow$  Model free

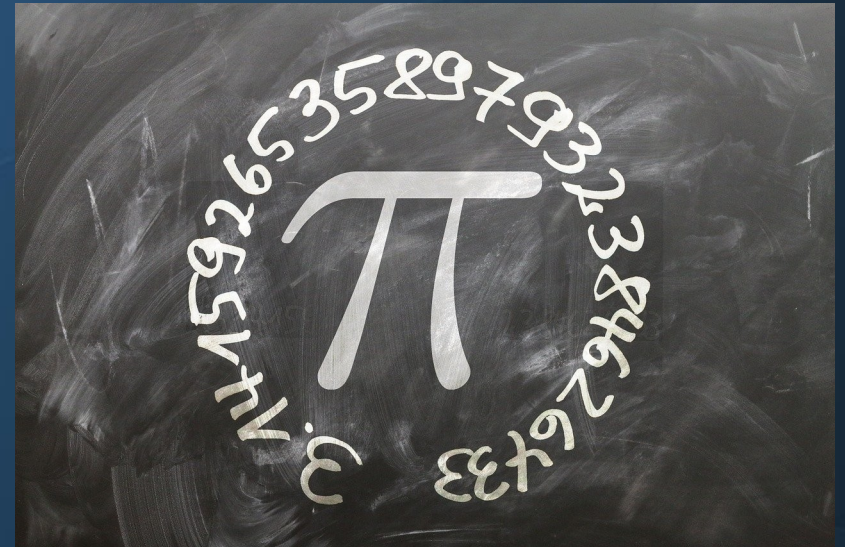


Exploration

# Explore Exploit Dilemma



Epsilon greedy (Q Learning)



Approximate Policy Directly

# On Policy vs. Off Policy

- One policy generates actions and updates value function → On policy
- One policy generates actions and another policy updates value function → Off policy
- Epsilon greedy → off policy learning
- Policy gradients → on policy learning

# Conclusions

- Keep track of rewards to estimate value and action value functions
- Recursive relationship between functions
- Have to interact w/ environment to learn dynamics
- Policy gradient & Actor critic → on policy model free



