

# Deep Reinforcement Learning Study Note

Han Cheol Moon  
School of Computer Science and Engineering  
Nanyang Technological University  
Singapore  
`hancheol001@edu.ntu.edu.sg`

March 5, 2023

# Contents

<b>1</b>	<b>Introduction</b>	<b>1</b>
1.1	Markov Chain . . . . .	1
1.2	Multi-Armed Bandits . . . . .	3
1.2.1	Action-value Methods . . . . .	3
1.2.2	Incremental Implementation of The Action-value Methods . . . . .	4
<b>2</b>	<b>Markov Decision Process</b>	<b>5</b>
2.1	Policies and Value Functions . . . . .	7
2.1.1	The State-Value Functions . . . . .	7
2.1.2	The Action-Value Function . . . . .	8
2.1.3	The Action-Advantage Function . . . . .	9
2.2	Optimality . . . . .	9
<b>3</b>	<b>Monte Carlo Method</b>	<b>12</b>
3.1	Intro . . . . .	12
3.2	First Visit vs. Every Visit . . . . .	12
3.3	On-Policy vs. Off-Policy . . . . .	12
3.4	On Policy . . . . .	13
3.5	Off Policy . . . . .	13
3.6	Temporal-Difference Learning . . . . .	14
3.6.1	Sarsa: On-policy TD Control . . . . .	14
3.6.2	Q-learning: Off-policy TD Control . . . . .	14
<b>4</b>	<b>Policy Gradient</b>	<b>15</b>

4.1	Natural Policy Gradient . . . . .	15
4.1.1	KL-divergence between perturbed distributions . . . . .	15
4.2	Proximal Policy Optimization . . . . .	17

Appendices

20

Appendix

20

A.1	Bellman Equation . . . . .	20
B.2	Importance Sampling . . . . .	22
C.3	Fisher Information . . . . .	23
D.4	Score Function . . . . .	23
E.5	Incremental Monte-Carlo . . . . .	24

# Chapter 1

## Introduction

Reinforcement learning (RL) is learning what to do so as to maximize a numerical reward signal. The learner is not told which actions to take, but instead must discover which actions yield the most reward by trying them. In the most interesting and challenging cases, actions may affect not only the immediate reward but also the next situation and, through that, all subsequent rewards.

Important features of RL:

- Trial-and-error search
- Delayed reward
- RL uses training information that *evaluates* the actions taken rather than *instructs* by giving correct actions.
  - Evaluative feedback, not about best or worst.

### 1.1 Markov Chain

- Reachable:  $i \rightarrow j$
- Communicate:  $i \leftrightarrow j$
- Irreducible:  $i \leftrightarrow j, \forall i, j$
- Absorbing state: If the only possible transition is to itself. This is also a terminal state.
- Transient state: A state  $s$  is called a transient state, if there is another state  $s'$ , that is reachable from  $s$ , but not vice versa.
- Recurrent state: A state that is not transient.
- Periodic: A state is periodic if all of the paths leaving  $s$  come back after some multiple steps ( $k > 1$ ).
  - Recurrent state is aperiodic if  $k = 1$ .
- Ergodicity if a Markov chain follows:

- Irreducible
- Recurrent
- Aperiodic

## 1.2 Multi-Armed Bandits

Bandit problems are stateless. Each arm has a fixed distribution of rewards. It does not depend on which arms were pulled previously. The goal is to explore the reward distributions of all arms and then keep pulling the best one.

Markov decision processes are a temporal extension of bandit problems: pulling an arm influences the future rewards. Technically, there is a state that changes by pulling an arm. The reward distributions depend on that state.

You can view bandit problems as Markov decision processes where all states are terminal. In that case, all decision sequences have a length of 1 and subsequent pulls don't influence each other.

- When the lever of a slot machine is pulled it gives a random reward coming from a probability distribution specific to that machine.
- Although the machines look identical, their reward probability distributions are different.
- In each turn, gamblers need to decide whether to play the machine that has given the highest average reward so far, or to try another machine.

In our  $k$ -armed bandit problem, each of the  $k$  actions has an expected or mean reward given that that action is selected: let us call this the *value* of that action. We denote the action selected on time step  $t$  as  $A_t$ , and the corresponding reward as  $R_t$ . The value then of an arbitrary action  $a$ , denoted  $q_*(a)$ , is the expected reward given that  $a$  is selected:

$$q_*(a) = \mathbb{E}[R_t | A_t = a].$$

Since we do not know which action is the best, we have to estimate the value of actions  $a$  at time step  $t$ ,  $Q_t(a)$ .

### 1.2.1 Action-value Methods

One natural way to estimate this is by averaging the rewards actually received:  $Q_t(a)$  is sum of rewards when  $a$  taken prior  $t$  over number of times  $a$  taken prior to  $t$ :

$$Q_t(a) = \frac{\sum_{i=1}^{t-1} R_i \cdot \mathbb{1}_{A_i=a}}{\sum_{i=1}^{t-1} \mathbb{1}_{A_i=a}}$$

Then, the simplest action selection rule is to select one of the actions with the highest estimated value, which is a greedy action selection method represented as follows:

$$A_t = \underset{a}{\operatorname{argmax}} Q_t(a).$$

This approach always exploits current knowledge to maximize immediate reward; it spends no time at all sampling apparently inferior actions to see if they might really be better. A simple alternative is to behave greedily most of the time, but every once in a while, instead select randomly from among all the actions. This near greedy action selection rule is called  $\epsilon$ -greedy method. In the limit as the number of steps increases, every action will be sampled an infinite number of times, which ensures all the  $Q_t(a)$  converge to  $q_*(a)$ .

### 1.2.2 Incremental Implementation of The Action-value Methods

$$\begin{aligned}
 Q_{n+1} &= \frac{1}{n} \sum_{i=1}^n R_i \\
 &= \frac{1}{n} \left( R_n + \sum_{i=1}^{n-1} R_i \right) \\
 &= \frac{1}{n} \left( R_n + (n-1) \frac{1}{n-1} \sum_{i=1}^{n-1} R_i \right) \\
 &= \frac{1}{n} \left( R_n + (n-1) Q_n \right) \\
 &= Q_n + \frac{1}{n} [R_n - Q_n]
 \end{aligned}$$

- $Q_n$ : Old estimate
- $Q_{n+1}$ : New estimate
- $R_n$ : New reward

This is an incremental formulas for updating averages with small, constant computation required to process each new reward. This update rule can be expressed in a general form:

$$\textit{NewEstimate} \leftarrow \textit{OldEstimate} + \textit{StepSize} \underbrace{\left[ \textit{Target} - \textit{OldEstimate} \right]}_{\textit{error}}.$$

The target is presumed to indicate a desirable direction in which to move, though it may be noisy.

## Chapter 2

# Markov Decision Process

The general framework of MDPs (representing environments as MDPs) allows us to model virtually any complex sequential decision-making problem under uncertainty in a way that RL agents can interact with and learn to solve solely through experience.

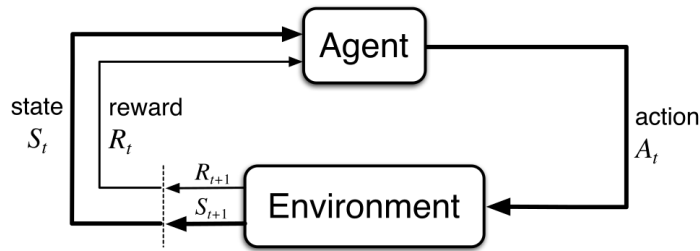


Figure 2.1: The agent-environment interaction in a Markov decision process.

Components of RL:

- An agent
- A policy
- A reward signal: what is good in an immediate sense.
- A value function: what is good in the long run.
- A model of the environment: This allows inferences to be made about how the environment will behave.

**Definition 1 (Markov Property)** A state  $S_t$  is *Markov* if and only if

$$P[S_{t+1}|S_t, A_t] = P[S_{t+1}|S_t, A_t, S_{t-1}, A_{t-1}, \dots]$$

- Actions: a mechanism to influence the environment
- State: specific configurations of the environment

**Definition 2 (Transition Function)**

$$p(s', r|s, a) = P(S_t = s', R_t = r|S_{t-1} = s, A_{t-1} = a)$$



- The way the environment changes as a response to actions is referred to as the state-transition probabilities, or more simply, the transition function, and is denoted by  $T(s, a, s')$ .
- $\sum_{r \in \mathcal{R}} \sum_{s' \in \mathcal{S}} p(s', r|s, a) = 1, \forall s \in \mathcal{S}, \forall a \in \mathcal{A}(s)$
- $p(s'|s, a) = P(S_t = s', R_t = r | S_{t-1} = s, A_{t-1} = a) = \sum_{r \in \mathcal{R}} p(s', r|s, a)$

**Definition 3 (Reward Function)**

$$\begin{aligned} r(s, a) &= \mathbb{E}[R_t | S_{t-1} = s, A_{t-1} = a] \\ &= \sum_{r \in \mathcal{R}} r \sum_{s' \in \mathcal{S}} p(s', r|s, a) \end{aligned}$$

- $R_t \in \mathcal{R} \subset \mathbb{R}$ . Note that negative reward is still reward.
- The expected reward function is defined as a function that takes in a state-action pair.
- It is the expectation of reward at time step  $t$ , given the state-action pair in the previous time step.
- It can also be defined as a function that takes a full transition tuple  $s, a, s'$ .

$$\begin{aligned} r(s, a, s') &= \mathbb{E}[R_t | S_{t-1} = s, A_{t-1} = a, S_t = s'] \\ &= \sum_{r \in \mathcal{R}} \frac{p(s', r|s, a)}{p(s'|s, a)} \end{aligned}$$

- **Reward Hypothesis:** That all of what we mean by goals and purposes can be well thought of as the maximization of the expected value of the cumulative sum of a received scalar signal (called reward).

**Definition 4 (Discount Factor,  $\gamma$ )**

$$G_t = R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \cdots + \gamma^{T-1} R_t$$

- The sum of all rewards obtained during the course of an episode is referred to as the return,  $G_t$ .
- Episodic tasks:  $G_t = R_{t+1} + R_{t+2} + \cdots + R_T$ .
  - $G_t = R_{t+1} + \gamma G_{t+1}$
- Continuing tasks:  $G_t = R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \cdots = \sum_{k=0}^{\infty} \gamma^k R_{t+k+1}$ .
  - $\gamma = 0$ : myopic evaluation
  - $\gamma = 1$ : far-sighted evaluation
  - Uncertainty about the future that may not be fully observed
  - Mathematically convenient to discount rewards.
  - Avoid infinite returns in cyclic Markov processes.

## 2.1 Policies and Value Functions

### 2.1.1 The State-Value Functions

Almost all reinforcement learning algorithms involve estimating *value functions*—functions of states (or of state-action pairs) that estimate how good it is for the agent to be in a given state (or how good it is to perform a given action in a given state). The notion of “how good” here is defined in terms of future rewards that can be expected, or, to be precise, in terms of **expected return**. Of course, the rewards the agent can expect to receive in the future depend on what actions it will take. Accordingly, value functions are defined with respect to particular ways of action, called policies.

Formally, a policy is a mapping from states to probabilities of selecting each possible action. It can be defined

$$\pi(a|s)$$

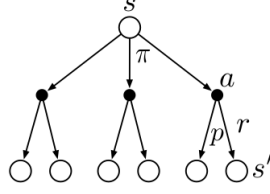
**Definition 5 (The State-Value Function,  $V$ )** *The state value function  $v(s)$  for policy  $\pi$  of an Markov Reward Process is the expected return starting from state  $s$*

$$v_{\pi}(s) = \mathbb{E}_{\pi}[G_t | S_t = s] = \mathbb{E}_{\pi} \left[ \sum_{k=0}^{\infty} \gamma^k R_{t+1+k} \middle| S_t = s \right], \quad \forall s \in \mathcal{S}$$

- The value of a state  $s$  is the expectation over policy  $\pi$ .
- Policies are universal plans, which provides all possible plans for all states.
  - Plans are not enough in stochastic environments.
  - Policy can be stochastic or deterministic.
  - A policy is a function that prescribes actions to take for a given non-terminal state.
- Reward: one-step signal that an agent gets/ Return: total discounted rewards/ Value function: expected return.
- If we are given a policy and the MDP, we should be able to calculate the expected return starting from every single state.
- *Bellman equation*

$$\begin{aligned} v_{\pi}(s) &= \mathbb{E}_{\pi}[G_t | S_t = s] \\ &= \mathbb{E}_{\pi}[R_{t+1} + \gamma G_{t+1} | S_t = s] \\ &= \mathbb{E}_{\pi}[R_{t+1} + \gamma v(s_{t+1}) | S_t = s] \\ &= \sum_a \pi(a|s) \sum_{s', r} p(s', r | s, a) [r + \gamma v_{\pi}(s')] \end{aligned}$$

- Starting from state  $s$ , the root node at the top, the agent could take any of some set of actions – three are shown in the diagram – based on its policy  $\pi$ . From each of these, the environment could respond with one of several next states,  $s'$  (two are shown in the figure), along with a reward,  $r$ , depending on its dynamics given by the function  $p$ . The Bellman equation averages over all the possibilities, weighting each by its probability of occurring. **It states that the value of the start state must equal (discounted) the value of the expected next state, plus the reward expected along the way.**



- The derivation details are given in Appendix A.1

### 2.1.2 The Action-Value Function

- Another critical question that we often need to ask is not merely about the value of a state but the value of taking action  $a$  at a state  $s$ .
- Which action is better under each policy?
- The action-value function, also known as  $Q$ -function or  $Q^\pi(s, a)$ , captures precisely this.
  - The expected return if the agent follows policy  $\pi$  after taking action  $a$  in state  $s$ .

**Definition 6 (The Action-Value Function,  $Q$ )** The action-value function  $q_\pi(s, a)$  for policy  $\pi$  is the expected return starting from state  $s$ , taking action  $a$  under policy  $\pi$

$$q_\pi(s, a) = \mathbb{E}_\pi[G_t | S_t = s, A_t = a] = \mathbb{E}_\pi \left[ \sum_{k=0}^{\infty} \gamma^k R_{t+k+1} \middle| S_t = s, A_t = a \right]$$

- The Bellman equation for action values is given by

$$\begin{aligned} q_\pi(s, a) &= \mathbb{E}_\pi[G_t | S_t = s, A_t = a] \\ &= \mathbb{E}_\pi[R_{t+1} + \gamma G_{t+1} | S_t = s, A_t = a] \\ &= \sum_{s', r} p(s', r | s, a) [r + \gamma v_\pi(s')] \end{aligned}$$

- Notice that we do not weight over actions because we are interested only in a specific action.
- The state-value function can be expressed by using the action-value function as follows:
- The derivation is given in Appendix A.1

$$\begin{aligned} v_\pi(s) &= \mathbb{E}_\pi[G_t | S_t = s] \\ &= \sum_{g_t} g_t \cdot P(G_t | S_t = s) \\ &= \sum_{g_t} g_t \sum_a P(G_t | S_t = s, A_t = a) P(A_t = a | S_t = s) \\ &= \sum_a \sum_{g_t} g_t P(G_t | S_t = s, A_t = a) P(A_t = a | S_t = s) \\ &= \sum_a q_\pi(s, a) \pi(a | s) \end{aligned}$$

Note that the expectation is parameterized  $\pi$  as written in  $\mathbb{E}_\pi$ . We can also prove it by the *Law of Total Expectation*,

$$\begin{aligned} v_\pi(s) &= \mathbb{E}_\pi[G_t | S_t = s] \\ &= \mathbb{E}[\mathbb{E}_\pi[G_t | S_t = s, A_t = a]] \\ &= \sum_a \mathbb{E}_\pi[G_t | S_t = s, A_t = a] P(A_t = a | S_t = s) \\ &= \sum_a q_\pi(s, a) \pi(a | s) \end{aligned}$$

An intuitive explanation of this derivation is that the expectation depends on an action  $a \sim \pi(a | s)$ . We want to estimate the expected total return by the sampled action (this is because the total return is the function of  $a$ , implicitly). Then we need to introduce an action variable  $a$  and its probability in the expression as the second line of the equation.

### 2.1.3 The Action-Advantage Function

**Definition 7 (The Action-Advantage Function,  $A$ )**

$$a_\pi(s, a) = q_\pi(s, a) - v_\pi(s)$$

- The advantage function describes how much better it is to take action  $a$  instead of following policy  $\pi$ . In other words, the advantage of choosing action  $a$  over the default action.

## 2.2 Optimality

Solving a reinforcement learning task means, roughly, finding a policy that achieves a lot of reward over the long run. For finite MDPs, we can precisely define an optimal policy in the following way. Value functions define a partial ordering over policies. A policy  $\pi$  is defined to be better than or equal to a policy  $\pi'$  if its expected return is greater than or equal to that of  $\pi'$  for all states. In other words,  $\pi \geq \pi'$  if and only if  $v_\pi(s) \geq v_{\pi'}(s)$  for all  $s \in \mathcal{S}$ . There is always at least one policy that is better than or equal to all other policies. This is an *optimal policy*. The optimal state-value function, denoted  $v_*$  can be defined as

**Definition 8 (Optimal State-Value Function)** *The optimal state-value function  $v_*(s)$  is the maximum value over all policies*

$$v_*(s) = \max_{\pi} v_\pi(s), \quad \forall s \in \mathcal{S}.$$

- The optimal state-value function can be obtained as follows:

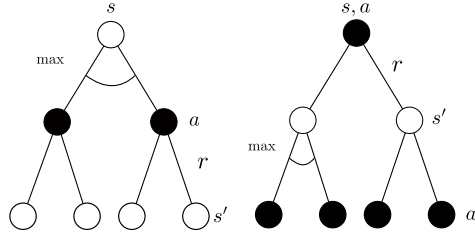
$$\begin{aligned}
v_*(s) &= \max_{a \in \mathcal{A}(s)} q_{\pi_*}(s, a) \\
&= \max_a \mathbb{E}_{\pi_*}[G_t | S_t = s, A_t = a] \\
&= \max_a \mathbb{E}_{\pi_*}[R_{t+1} + \gamma G_{t+1} | S_t = s, A_t = a] \\
&= \max_a \mathbb{E}_{\pi_*}[R_{t+1} + \gamma v_*(S_{t+1}) | S_t = s, A_t = a] \\
&= \max_a \sum_{r, s'} p(s', r | s, a) [r + \gamma v_*(s')]
\end{aligned}$$

**Definition 9 (Optimal Action-Value Function)** The optimal action-value function  $q_*(s, a)$  is the maximum value over all policies

$$q_*(s, a) = \max_{\pi} q_{\pi}(s, a), \quad \forall s \in \mathcal{S} \text{ and } a \in \mathcal{A}.$$

- The optimal action-value function can be obtained as follows:

$$\begin{aligned}
q_*(s, a) &= \mathbb{E}[R_{t+1} + \gamma v_*(S_{t+1}) | S_t = s, A_t = a] \\
&= \sum_{s', r} p(s', r | s, a) [r + \gamma \max_{a'} q_*(s', a')]
\end{aligned}$$



- The optimal value function specifies the best possible performance in the MDP.
- The MDP is solved when we know the optimal value function

**Theorem 1 (Optimal Policy Theorem)**

$$\pi \geq \pi' \quad \text{if} \quad v_{\pi}(s) \geq v_{\pi'}(s), \forall s$$

For any Markov Decision Process:

- There exists an optimal policy  $\pi_*$  that is better than or equal to all other policies,  $\pi_* \geq \pi, \forall \pi$
- All optimal policies achieve the optimal value function,  $v_{\pi_*}(s) = v_*(s)$
- All optimal policies achieve the optimal action-value function,  $q_{\pi_*}(s, a) = q_*(s, a)$

An optimal policy can be found by maximizing over  $q_*(s, a)$ ,

$$\pi_*(a|s) = \begin{cases} 1 & \text{if } a = \operatorname{argmax}_{a \in \mathcal{A}} q_*(s, a) \\ 0 & \text{otherwise} \end{cases}$$

- There is always a deterministic optimal policy for any MDP
- If we know  $q_*(s, a)$ , we immediately have the optimal policy
  - Q-learning: learns Q values first
  - Policy gradient: learns optimal policy without learning Q values

## Chapter 3

# Monte Carlo Method

### 3.1 Intro

- Prediction problem: refers to the problem of evaluating policies, of estimating value functions given a policy (learning to predict returns).
- Control problem: problem of finding optimal policies. Usually solved by the pattern of generalized policy iteration (GPI), where the competing processes of policy evaluation and policy improvement progressively move policies towards optimality.
- Policy evaluation: refers to algorithms that solve the prediction problem.
- Policy improvement: algorithms that make new policies that improve on an original policy by making it greedier than the original with respect to the value function of that original policy.

### 3.2 First Visit vs. Every Visit

The first-visit and the every-visit Monte-Carlo (MC) algorithms are both used to solve the prediction problem (or, also called, “evaluation problem”), that is, the problem of estimating the **value function** associated with a given (as input to the algorithms) fixed (that is, it does not change during the execution of the algorithm) policy, denoted by  $\pi$ . In general, even if we are given the policy  $\pi$ , we are not necessarily able to find the exact corresponding value function, so these two algorithms are used to estimate the value function associated with  $\pi$ .

Intuitively, we care about the value function associated with  $\pi$  because we might want or need to know “how good it is to be in a certain state”, if the agent behaves in the environment according to the policy  $\pi$ .

### 3.3 On-Policy vs. Off-Policy

The more episodes are collected, the better because the estimates of the functions will be. However, there is a problem. If the algorithm for policy improvement always updates the policy greedily, meaning it takes only actions leading to immediate reward, actions and states not on the

greedy path will not be sampled sufficiently, and potentially better rewards would stay hidden from the learning process.

Essentially, we are forced to make a choice between making the best decision given the current information or start exploring and finding more information. This is also known as the Exploration vs. Exploitation Dilemma.

We are looking for something like a middle ground between those. Full-on exploration would mean that we would need a lot of time to collect the needed information, and full-on exploitation would make the agent stuck into a local reward maximum. There are two approaches to ensure all actions are sampled sufficiently called on-policy and off-policy methods.

### 3.4 On Policy

On-policy methods solve the exploration vs exploitation dilemma by including randomness in the form of a policy that is soft, meaning that non-greedy actions are selected with some probability. These policies are called  $\epsilon$ -greedy policies as they select random actions with an  $\epsilon$  probability and follow the optimal action with  $1 - \epsilon$  probability

Since the probability of selecting from the action space randomly is  $\epsilon$ , the probability of selecting any particular non-optimal (non-greedy) action is  $\epsilon/|\mathcal{A}(s)|$ . The probability of following the optimal action will always be slightly higher, however, because we have a  $1 - \epsilon$  probability of selecting it outright and  $\epsilon/|\mathcal{A}(s)|$  probability of selecting it from sampling the action space <sup>1</sup>:

$$P(a_t^*) = 1 - \epsilon + \epsilon/|\mathcal{A}(s)|.$$

It is also worth noting that because the optimal action will be sampled more often than the others making on-policy algorithms will generally converge faster but they also have the risk of trapping the agent into a local optimum of the function.

### 3.5 Off Policy

All learning control methods face a dilemma: They seek to learn action values conditional on subsequent *optimal* behavior, but they need to behave non-optimally in order to explore all actions (to find the optimal actions). How can they learn about the optimal policy while behaving according to an exploratory policy? The on-policy approach in the preceding section is actually a compromise. It learns action values not for the optimal policy, but for a near-optimal policy that still explores. A more straightforward approach is to use two policies, one that is learned about and that becomes the optimal policy, and one that is more exploratory and is used to generate behavior. The policy being learned about is called the *target policy*, and the policy used to generate behavior is called the *behavior policy*. In this case we say that learning is from data “off” the target policy, and the overall process is termed *off-policy learning*.

Given a starting state  $S_t$ , the probability of the subsequent state-action trajectory,  $A_t, S_{t+1}, A_{t+1}, \dots, S_T$ ,

---

<sup>1</sup>Since the greedy (or optimal) action can be selected by either  $1-\epsilon$  or  $\epsilon/|\mathcal{A}(s)|$



occurring under any policy  $\pi$  is

$$\begin{aligned} P(A_t, S_{t+1}, A_{t+1}, \dots, S_T) &= \pi(A_t|S_t)p(S_{t+1}|S_t, A_t)\pi(A_{t+1}|S_{t+1})(S_T|S_{T-1}, A_{T-1}) \\ &= \prod_{k=t}^{T-1} \pi(A_k|S_k)p(S_{k+1}|S_k, A_k), \end{aligned}$$

where  $p$  here is the state-transition probability function. Thus, the relative probability of the trajectory under the target and behavior policies is

$$\begin{aligned} \rho_{t:T-1} &= \frac{\prod_{k=t}^{T-1} \pi(A_k|S_k)p(S_{k+1}|S_k, A_k)}{\prod_{k=t}^{T-1} b(A_k|S_k)p(S_{k+1}|S_k, A_k)} \\ &= \prod_{k=t}^{T-1} \frac{\pi(A_k|S_k)}{b(A_k|S_k)}. \end{aligned}$$

The importance sampling ratio ends up depending only on the two policies and the sequence, not on the MDP (state-transition probability).

$$\mathbb{E}[G_t|S_t = s] = v_b(s).$$

Recall that we wish to

$$\mathbb{E}[\rho_{t:T-1}G_t|S_t = s] = v_\pi(s).$$

## 3.6 Temporal-Difference Learning

TD learning is a combination of Monte-Carlo method and dynamic programming ideas. It learns directly from raw experience without a model of the environment's dynamics.

Constant- $\alpha$  MC:

$$V(S_t) \leftarrow V(S_t) + \alpha [G_t - V(S_t)]$$

One-step TD (or TD(0)):

$$V(S_t) \leftarrow V(S_t) + \alpha \left[ \underbrace{R_{t+1} + \gamma V(S_{t+1}) - V(S_t)}_{\text{TD error}} \right]$$

### 3.6.1 Sarsa: On-policy TD Control

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha [R_{t+1} + \gamma Q(S_{t+1}, A_{t+1}) - Q(S_t, A_t)]$$

### 3.6.2 Q-learning: Off-policy TD Control

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha [R_{t+1} + \gamma \max_a Q(S_{t+1}, a) - Q(S_t, A_t)]$$

## Chapter 4

# Policy Gradient

We will use a gradient ascent algorithm:

$$\begin{aligned} J(\theta) &= \mathbb{E}_{\tau \sim \pi_{\theta}(\tau)}[r(\tau)] \\ &= \int \pi_{\theta}(\tau) r(\tau) \end{aligned}$$

It is a expected reward under the policy  $\pi_{\theta}$ .

$$\theta \leftarrow \theta + \eta \nabla_{\theta} J(\theta)$$

Note that by using REINFORCE algorithm which can be expressed as follows:

$$\nabla_{\theta} \pi_{\theta}(\tau) = \pi_{\theta} \frac{\nabla_{\theta} \pi_{\theta}(\tau)}{\pi_{\theta}(\tau)} = \pi_{\theta} \nabla_{\theta} \log \pi_{\theta}(\tau)$$

We can express  $J(\theta)$  as follows:

$$\begin{aligned} \nabla_{\theta} J(\theta) &= \int \nabla_{\theta} \pi_{\theta}(\tau) r(\tau) = \int \pi_{\theta} \nabla_{\theta} \log \pi_{\theta}(\tau) r(\tau) \\ &= \mathbb{E}_{\tau \sim \pi_{\theta}(\tau)}[\nabla_{\theta} \log \pi_{\theta}(\tau) r(\tau)] \end{aligned}$$

Note that  $\nabla_{\theta} \log \pi_{\theta}(\tau)$  is the maximum loglikelihood of trajectory, because gradient is the maximum direction of the function. This expectation can be estimated by Monte-Carlo method.

## 4.1 Natural Policy Gradient

### 4.1.1 KL-divergence between perturbed distributions

Let  $p(x; \theta)$  be some family of probability distributions over  $x$  parameterized by a vector of real numbers  $\theta$ . We're interested in knowing how much the distribution changes when we perturb the parameter vector from a fixed  $\theta_t$  to some new value  $\theta_t + \delta\theta$ . As a measure of change in probability distribution, we can use the KL-divergence measure. Specifically, we want to measure  $D_{KL}(p(x; \theta_t) || p(x; \theta_t + \delta\theta))$ , but we want to write it in a form amenable to the gradient-based

update formulation. We can do this by taking it's second-order Taylor expansion around  $\theta_t$ . During the derivation, we'll find that a lot of terms in the expansion disappear leaving us with a very simple expression that's perfect for our purposes.

Looking first at the full KL-divergence, we see that the term we want to

$$\begin{aligned} D_{KL}(p(x; \theta_t) || p(x; \theta_t + \delta\theta)) &= \int p(x; \theta_t) \log \frac{p(x; \theta_t)}{p(x; \theta_t + \delta\theta)} dx \\ &= \int p(x; \theta_t) \log p(x; \theta_t) dx - \int p(x; \theta_t) \log p(x; \theta_t + \delta\theta) dx \end{aligned}$$

Note that the second-order Taylor series expansion is

$$f(\theta) \approx f(\theta_t) + \nabla f(\theta_t)^T \delta\theta + \frac{1}{2} \delta\theta^T \nabla^2 f(\theta_t) \delta\theta,$$

where  $\theta = \theta_t + \delta\theta$ , or equivalently  $\delta\theta = \theta - \theta_t$ . Applying that expansion to the pertinent term in the  $KL$ -divergence expression, we get

$$\log p(x; \theta_t + \delta\theta) \approx \log p(x; \theta_t) + \left( \frac{\nabla p(x; \theta_t)}{p(x; \theta_t)} \right)^T \theta + \frac{1}{2} \delta\theta^T (\nabla^2 \log p(x; \theta_t)) \delta\theta.$$

Pugging this second-order Taylor expansion back into the above expression for the  $D_{KL}$  gives

$$\begin{aligned} D_{KL}(p(x; \theta_t) || p(x; \theta_t + \delta\theta)) &= \int p(x; \theta_t) \log p(x; \theta_t) dx - \int p(x; \theta_t) \log p(x; \theta_t + \delta\theta) dx \\ &\approx \int p(x; \theta_t) \log p(x; \theta_t) dx \\ &\quad - \int p(x; \theta_t) \left( \log p(x; \theta_t) + \left( \frac{\nabla p(x; \theta_t)}{p(x; \theta_t)} \right)^T \theta + \frac{1}{2} \delta\theta^T (\nabla^2 \log p(x; \theta_t)) \delta\theta \right) dx \\ &= 0 - \underbrace{\int p(x; \theta_t) \left( \frac{\nabla p(x; \theta_t)}{p(x; \theta_t)} \right)^T dx}_{=0} \theta - \frac{1}{2} \delta\theta^T \left( \int p(x; \theta_t) \nabla^2 \log p(x; \theta_t) dx \right) \delta\theta \\ &= -\frac{1}{2} \delta\theta^T \left( \int p(x; \theta_t) \nabla^2 \log p(x; \theta_t) dx \right) \delta\theta. \end{aligned}$$

The Hessian can be computed as follows:

$$\begin{aligned} \frac{\partial^2}{\partial \theta_t^i \partial \theta_t^j} [\log p(x; \theta_t)] &= \frac{\partial}{\partial \theta_t^i} \left( \frac{\frac{\partial}{\partial \theta_t^j} p(x; \theta_t)}{p(x; \theta_t)} \right) \\ &= \frac{p(x; \theta_t) \frac{\partial^2}{\partial \theta_t^i \partial \theta_t^j} p(x; \theta_t) - \frac{\partial}{\partial \theta_t^i} p(x; \theta_t) \frac{\partial}{\partial \theta_t^j} p(x; \theta_t)}{p(x; \theta_t)^2} \\ &= \frac{1}{p(x; \theta_t)} \frac{\partial^2}{\partial \theta_t^i \partial \theta_t^j} p(x; \theta_t) - \left( \frac{\frac{\partial}{\partial \theta_t^i} p(x; \theta_t)}{p(x; \theta_t)} \right) \left( \frac{\frac{\partial}{\partial \theta_t^j} p(x; \theta_t)}{p(x; \theta_t)} \right). \end{aligned}$$

The second term is an element of the outer product between  $\nabla \log p(x; \theta_t)$  and itself. In matrix form, this becomes

$$\nabla^2 \log p(x; \theta_t) = \frac{1}{p(x; \theta_t)} \nabla^2 p(x; \theta_t) - \nabla \log p(x; \theta_t) \nabla \log p(x; \theta_t)^T.$$

Finally, we get

$$\begin{aligned} D_{KL}(p(x; \theta_t) || p(x; \theta_t + \delta\theta)) &\approx -\frac{1}{2} \delta\theta^T \int p(x; \theta_t) \nabla^2 \log p(x; \theta_t) dx \delta\theta \\ &= \frac{1}{2} \delta\theta^T \left( \int \nabla^2 \log p(x; \theta_t) dx \right) \delta\theta \\ &\quad + \frac{1}{2} \delta\theta^T \left( \int p(x; \theta_t) [\nabla \log p(x; \theta_t) \nabla \log p(x; \theta_t)^T] dx \right) \delta\theta \\ &= \frac{1}{2} \delta\theta^T \underbrace{\left( \int p(x; \theta_t) [\nabla \log p(x; \theta_t) \nabla \log p(x; \theta_t)^T] dx \right)}_{G(\theta_t)} \delta\theta. \end{aligned}$$

The central matrix here  $G(\theta_t)$  is known as the **Fisher Information matrix** and can has been thoroughly studied within the field of Information Geometry as the natural Riemannian structure on a manifold of probability distributions. As such it defines a natural norm on perturbations to probability distributions, which was our original motivation for examining the second-order Taylor expansion of the KL-divergence in the first place.

$$\theta_{t+1} = \theta_t - \eta_t G(\theta_t)^{-1} \nabla f(\theta_t).$$

## 4.2 Proximal Policy Optimization

PPO objective is

$$\theta_{k+1} = \operatorname{argmax}_{\theta} E_{s,a \sim \pi_{\theta_k}} [L(s, a, \theta_k, \theta)],$$

where  $L$  is given by

$$L(s, a, \theta_k, \theta) = \min \left( \frac{\pi_{\theta}(a|s)}{\pi_{\theta_k}(a|s)} A^{\pi_{\theta_k}}(s, a), \operatorname{Clip} \left( \frac{\pi_{\theta}(a|s)}{\pi_{\theta_k}(a|s)}, 1 - \varepsilon, 1 + \varepsilon \right) A^{\pi_{\theta_k}}(s, a) \right).$$

Roughly,  $\varepsilon$  is a hyperparameter which says how far away the new policy is allowed to go from the old one. A simpler expression of the above expression is

$$L(s, a, \theta_k, \theta) = \min \left( \frac{\pi_{\theta}(a|s)}{\pi_{\theta_k}(a|s)} A^{\pi_{\theta_k}}(s, a), g(\varepsilon, A^{\pi_{\theta_k}}(s, a)) \right), \quad (4.1)$$

where

$$g(\varepsilon, A) = \begin{cases} (1 + \varepsilon)A & A \geq 0 \\ (1 - \varepsilon)A & A < 0. \end{cases} \quad (4.2)$$

**Positive Advantage:** Suppose the advantage for that state-action pair is positive, in which case its contribution to the objective reduces to

$$L(s, a, \theta_k, \theta) = \min \left( \frac{\pi_\theta(a|s)}{\pi_{\theta_k}(a|s)}, 1 + \varepsilon \right) A^{\pi_{\theta_k}}(s, a). \quad (4.3)$$

Because the advantage is positive, the objective will increase if the action becomes more likely—that is, if  $\pi_\theta(a|s)$  increases. But the min in this term puts a limit to how much the objective can increase. Once  $\pi_\theta(a|s) > (1 + \varepsilon)\pi_{\theta_k}(a|s)$ , the min kicks in and this term hits a ceiling of  $(1 + \varepsilon)A^{\pi_{\theta_k}}(s, a)$ . Thus, the new policy does not benefit by going far away from the old policy.

**Negative Advantage:** Suppose the advantage for that state-action pair is negative, in which case its contribution to the objective reduces to

$$L(s, a, \theta_k, \theta) = \min \left( \frac{\pi_\theta(a|s)}{\pi_{\theta_k}(a|s)}, 1 - \varepsilon \right) A^{\pi_{\theta_k}}(s, a). \quad (4.4)$$

Because the advantage is negative, the objective will increase if the action becomes less likely—that is, if  $\pi_\theta(a|s)$  decreases. But the max in this term puts a limit to how much the objective can increase. Once  $\pi_\theta(a|s) < (1 - \varepsilon)\pi_{\theta_k}(a|s)$ , the max kicks in and this term hits a ceiling of  $(1 - \varepsilon)A^{\pi_{\theta_k}}(s, a)$ . Thus, again, the new policy does not benefit by going far away from the old policy.

What we have seen so far is that clipping serves as a regularizer by removing incentives for the policy to change dramatically, and the hyperparameter  $\varepsilon$  corresponds to how far away the new policy can go from the old while still profiting the objective.

# Bibliography

# Appendix

## A.1 Bellman Equation

*Bellman equation* can be derived as follows:

$$\begin{aligned}
v_\pi(s) &= \mathbb{E}_\pi[G_t | S_t = s] \\
&= \mathbb{E}_\pi[R_{t+1} + \gamma G_{t+1} | S_t = s] \\
&= \mathbb{E}_\pi[R_{t+1} | S_t = s] + \mathbb{E}_\pi[\gamma G_{t+1} | S_t = s], \quad \text{By Linearity of Expectation.} \\
&= \sum_{r_{t+1}} r_{t+1} P(R_{t+1} | S_t = s) + \mathbb{E}_\pi[\gamma G_{t+1} | S_t = s] \\
&= \sum_{r_{t+1}} r_{t+1} \sum_a P(R_{t+1} | S_t = s, A_t = a) P(A_t = a | S_t = s) + \mathbb{E}_\pi[\gamma G_{t+1} | S_t = s] \\
&= \sum_a \sum_{r_{t+1}} r_{t+1} \sum_{s'} P(R_{t+1}, S_{t+1} = s' | S_t = s, A_t = a) P(A_t = a | S_t = s) + \mathbb{E}_\pi[\gamma G_{t+1} | S_t = s] \\
&= \sum_a \sum_r r \sum_{s'} P(s', r | s, a) \pi(a | s) + \mathbb{E}_\pi[\gamma G_{t+1} | S_t = s] \\
&= \sum_a \sum_{s'} \sum_r r P(s', r | s, a) \pi(a | s) + \mathbb{E}_\pi[\gamma G_{t+1} | S_t = s] \\
&= \sum_a \sum_{s'} \mathcal{P}_{ss'}^a \mathcal{R}_{ss'}^a \pi(a | s) + \gamma \mathbb{E}_\pi[G_{t+1} | S_t = s] \\
&= \sum_a \sum_{s'} \mathcal{P}_{ss'}^a \mathcal{R}_{ss'}^a \pi(a | s) + \gamma \sum_a \mathbb{E}_\pi[G_{t+1} | S_t = s, A_t = a] P(A_t | S_t) \\
&= \sum_a \pi(a | s) \left[ \sum_{s'} \mathcal{P}_{ss'}^a \mathcal{R}_{ss'}^a + \gamma \mathbb{E}_\pi[G_{t+1} | S_t = s, A_t = a] \right] \\
&= \sum_a \pi(a | s) \left[ \sum_{s'} \mathcal{P}_{ss'}^a \mathcal{R}_{ss'}^a + \gamma \sum_{g_{t+1}} g_{t+1} P(G_{t+1} | S_t = s, A_t = a) \right] \\
&= \sum_a \pi(a | s) \left[ \sum_{s'} \mathcal{P}_{ss'}^a \mathcal{R}_{ss'}^a + \gamma \sum_{g_{t+1}} g_{t+1} \frac{P(G_{t+1}, S_t = s, A_t = a)}{P(S_t = s, A_t = a)} \right] \\
&= \sum_a \pi(a | s) \left[ \sum_{s'} \mathcal{P}_{ss'}^a \mathcal{R}_{ss'}^a + \gamma \sum_{g_{t+1}} g_{t+1} \frac{\sum_{s'} P(G_{t+1}, S_t = s, S_{t+1} = s', A_t = a)}{P(S_t = s, A_t = a)} \right] \\
&= \sum_a \pi(a | s) \left[ \sum_{s'} \mathcal{P}_{ss'}^a \mathcal{R}_{ss'}^a + \gamma \sum_{g_{t+1}} g_{t+1} \frac{\sum_{s'} P(G_{t+1} | s, s', a) P(s, s', a)}{P(s, a)} \right] \\
&= \sum_a \pi(a | s) \left[ \sum_{s'} \mathcal{P}_{ss'}^a \mathcal{R}_{ss'}^a + \gamma \sum_{g_{t+1}} g_{t+1} \sum_{s'} P(G_{t+1} | s, s', a) P(s' | s, a) \right] \\
&= \sum_a \pi(a | s) \left[ \sum_{s'} \mathcal{P}_{ss'}^a \mathcal{R}_{ss'}^a + \gamma \sum_{s'} P(s' | s, a) \sum_{g_{t+1}} g_{t+1} P(G_{t+1} | s') \right] \quad \text{by Markov Property} \\
&= \sum_a \pi(a | s) \left[ \sum_{s'} \mathcal{P}_{ss'}^a \mathcal{R}_{ss'}^a + \gamma \sum_{s'} \mathcal{P}_{ss'}^a v_\pi(s') \right] \\
&= \sum_a \pi(a | s) \sum_{s'} \mathcal{P}_{ss'}^a \left[ \mathcal{R}_{ss'}^a + \gamma v_\pi(s') \right] \\
&= \sum_a \pi(a | s) \sum_{r, s'} p(s', r | s, a) \left[ r + \gamma v_\pi(s') \right]
\end{aligned}$$



Or simply,

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

Reference

Similarly,

$$\begin{aligned} q_\pi(s, a) &= \mathbb{E}_\pi[G_t | S_t = s, A_t = a] \\ &= \mathbb{E}_\pi[R_{t+1} + \gamma G_{t+1} | S_t = s, A_t = a] \\ &= \mathbb{E}_\pi[R_{t+1} | S_t = s, A_t = a] + \gamma \mathbb{E}_\pi[G_{t+1} | S_t = s, A_t = a] \\ &= \sum_r r p(r|s, a) + \gamma \mathbb{E}_\pi[G_{t+1} | S_t = s, A_t = a] \\ &= \sum_r r \sum_{s'} p(s', r|s, a) + \gamma \mathbb{E}_\pi[G_{t+1} | S_t = s, A_t = a] \\ &= \sum_{s', r} r p(s', r|s, a) + \gamma \mathbb{E}[\mathbb{E}_\pi[G_{t+1} | S_t = s, A_t = a, R_{t+1}, S_{t+1}]] \quad \text{By Law of Total Expectation.} \\ &= \sum_{s', r} r p(s', r|s, a) + \gamma \sum_{s', r} \mathbb{E}_\pi[G_{t+1} | S_t = s, A_t = a, R_{t+1} = r, S_{t+1} = s'] p(s', r|s, a) \\ &= \sum_{s', r} p(s', r|s, a) [r + \gamma \mathbb{E}_\pi[G_{t+1} | S_t = s, A_t = a, R_{t+1} = r, S_{t+1} = s']] \\ &= \sum_{s', r} p(s', r|s, a) [r + \gamma \mathbb{E}_\pi[G_{t+1} | S_{t+1} = s']] \quad \text{By Markov Property.} \\ &= \sum_{s', r} p(s', r|s, a) [r + \gamma v_\pi(s')] \end{aligned}$$

## B.2 Importance Sampling

You have two distributions,  $P(A)$  and  $P(B)$ , and you have a sequence sampled from  $A$ . You can estimate an expectation of  $A$ ,

$$E[A] = \sum p(a) h(a).$$

Can we use the above equation for estimating an expectation of  $B$ ? Yes.

$$E[B] = \sum \frac{p(b)}{p(a)} h(a).$$

The ratio  $\frac{p(b)}{p(a)}$  tells us how likely to observe some results under  $p(b)$  compared to  $p(a)$ .

### C.3 Fisher Information

Suppose we have a model parameterized by parameter vector  $\theta$  that models a distribution  $p(x; \theta)$ . In frequentist statistics, the way we learn  $\theta$  is to maximize the likelihood of  $p(x; \theta)$ . To assess the goodness of our estimate of  $\theta$  we define a **score function** as follows:

$$f(\theta) = \nabla_{\theta} \log p(x; \theta).$$

The expected value of score function is zero.

$$\begin{aligned} \mathbb{E}_{p(x; \theta)}[f(\theta)] &= \mathbb{E}_{p(x; \theta)}[\nabla_{\theta} \log p(x; \theta)] \\ &= \int p(x; \theta) \nabla_{\theta} \log p(x; \theta) dx \\ &= \int p(x; \theta) \frac{\nabla_{\theta} p(x; \theta)}{p(x; \theta)} dx \\ &= 0 \end{aligned}$$

The covariance of the score function is given by

$$\text{Cov}[f(\theta), f(\theta)] = \mathbb{E}_{p(x; \theta)}[(f(\theta) - 0)(f(\theta) - 0)^T] = \text{Var}[f(\theta), f(\theta)].$$

This the definition of Fisher information and it can be written

$$F = \mathbb{E}_{p(x; \theta)}[\nabla \log p(x; \theta) \nabla \log p(x; \theta)^T].$$

Empirically,

$$F = \frac{1}{N} \sum_{i=1}^N \nabla \log p(x; \theta) \nabla \log p(x; \theta)^T.$$

### D.4 Score Function

In statistics, *the score (or informant) is the gradient of the log-likelihood function with respect to the parameter vector.*

- Evaluated at a particular point of the parameter vector, the score indicates the **steepness of the log-likelihood function and thereby the sensitivity to infinitesimal changes to the parameter values.**
- If the log-likelihood function is continuous over the parameter space, the score will **vanish at a local maximum or minimum**; this fact is used in maximum likelihood estimation to find the parameter values that maximize the likelihood function.

Since the score is a function of the observations that are subject to sampling error, it lends itself to a test statistic known as score test in which the parameter is held at a particular value. Further, the ratio of two likelihood functions evaluated at two distinct parameter values can be understood as a definite integral of the score function.[2]

## E.5 Incremental Monte-Carlo

Incremental Mean:

$$\mu_k = \frac{1}{k} \sum_{j=1}^k x_j \quad (5)$$

$$= \frac{1}{k} \left( x_k + (k-1) \frac{1}{(k-1)} \sum_{j=1}^{k-1} x_j \right) \quad (6)$$

$$= \frac{1}{k} \left( x_k + (k-1) \mu_{k-1} \right) \quad (7)$$

$$= \mu_{k-1} + \frac{1}{k} (x_k - \mu_{k-1}) \quad (8)$$

Incremental MC:

- $N_{n+1}(S_t^n) = N_n(S_t^n) + 1$  for rest  $N_{n+1}(s) = N_n(s)$
- $V_{n+1}(S_t) = V_n(S_t) + \frac{G_{t:T}^n - V_n(S_t)}{N_n(S_t)}$  for rest  $V_{n+1}(s) = V_n(s)$
- $V_{n+1}(S_t) = V_n(S_t) + \alpha(G_{t:T}^n - V_n(S_t))$  for rest  $V_{n+1}(s) = V_n(s)$