



Reinforcement Learning Notes

Learn from trying!

Author: occupymars

Date: June. 20, 2025

Version: 0.1

Contents

Chapter 1 Math of Reinforcement Learning	1
1.1 Notation	1
1.2 Markov Decision Process	1
1.3 Value Function	1
1.4 Solving Value Function	2
1.5 Action Value Function	3
1.6 Bellman Optimality Equation	3
1.7 Value Iteration and Policy Iteration	6
1.8 Monte Carlo Method	8
1.9 Stochastic Approximation	10
1.10 Temporal-Difference Methods	12
1.11 Value Function Methods	16
1.12 Policy Gradient Methods	20
1.13 Actor-Critic Methods	25
1.14 Appendix: Measure-Theoretic Probability Theory	30
Chapter 2 From LQR to RL	33
2.1 LQR and Value function	33

Chapter 1 Math of Reinforcement Learning

Introduction

- Markov Decision Process
- Value Function
- Solving Value Function

- Action Value Function
- Bellman Optimality Equation

1.1 Notation

At some point I use V^π to denote the value function under a certain policy π , but later I use v_π , they are the same. It's a lagacy from my reading of the *Dive into Deep Learning* book at first.

1.2 Markov Decision Process

State and **Action** can describe a robot state respect to the enviroment and actions to move around, \mathcal{S}, \mathcal{A} are states and actions a robot can take, when taking an action, state after may not be deterministic, it has a probability. We use a transition function $T : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow [0, 1]$ to denote this, $T(s, a, s') = p(s' | s, a)$ is the probability of reaching s' given s and a . For $\forall s \in \mathcal{S}$ and $\forall a \in \mathcal{A}$, $\sum_{s' \in \mathcal{S}} T(s, a, s') = 1$.

Reward $r : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$, $r(s, a)$ depends on current state and action. And the reward may also be stochastic, given state and action, the reward has probability $p(r | s, a)$.

Policy $\pi(a | s)$ tells agent which actions to take at every state, $\sum_a \pi(a | s) = 1$.

This can build a Markov Decision Process, $(\mathcal{S}, \mathcal{A}, \mathcal{T}, r)$ from the **Trajectory** $\tau = (s_0, a_0, r_0, s_1, a_1, r_1, s_2, a_2, r_2, \dots)$, which has probability of:

$$p(\tau) = \pi(a_0 | s_0) \cdot p(s_1 | s_0, a_0) \cdot \pi(a_1 | s_1) \cdot p(s_2 | s_1, a_1) \cdots$$

We then define **Return** as the total reward $R(\tau) = \sum_t r_t$, the goal of reinforcement learning is to find a trajectory that has the largest return. The trajectory might be infinite, so in order for a meaningful formular of its return, we introduce a discount factor $\gamma < 1$, $R(\tau) = \sum_{t=0}^{\infty} \gamma^t r_t$. For large γ , the robot is encouraged to explore, for small one to take a short trajectory to goal.

Markov system only depend on current state and action, not the history one (but we can always augment the system).

Remark In this book, the reward to a state that is not forbidden and is not a boundary is set to 0. There are five actions: up, right, down, left, and stay, with a_1, \dots, a_5 .

1.3 Value Function

Value Function is the value of a state, from that state, the expected sum reward (return).

The formular of value function is:

$$V^\pi(s_0) = \mathbb{E}_{a_t \sim \pi(s_t)}[R(\tau)] = \mathbb{E}_{a_t \sim \pi(s_t)} \left[\sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) \right] \quad (1.1)$$

If we divede the trajectory into two parts, s_0 and τ' , we get the return:

$$R(\tau) = r(s_0, a_0) + \gamma \sum_{t=1}^{\infty} \gamma^{t-1} r(s_t, a_t) = r(s_0, a_0) + \gamma R(\tau')$$

Put it back into the value function, using law of total expectation:

$$\mathbb{E}[X] = \sum_a \mathbb{E}[X | A = a] p(a) = \mathbb{E}_a [\mathbb{E}[X | A = a]]$$

we get:

$$\begin{aligned}
V^\pi(s_0) &= \mathbb{E}_{a_t \sim \pi(s_t)} [r(s_0, a_0) + \gamma R(\tau')] \\
&= \mathbb{E}_{a_0 \sim \pi(s_0)} [r(s_0, a_0)] + \gamma \mathbb{E}_{a_t \sim \pi(s_t)} [R(\tau')] \\
&= \mathbb{E}_{a_0 \sim \pi(s_0)} [r(s_0, a_0)] + \gamma \mathbb{E}_{a_0 \sim \pi(s_0)} [\mathbb{E}_{s_1 \sim p(s_1|a_0, s_0)} [\mathbb{E}_{a_t \sim \pi(s_t)} [R(\tau') \mid s_1, a_0]]] \\
&= \mathbb{E}_{a_0 \sim \pi(s_0)} [r(s_0, a_0)] + \gamma \mathbb{E}_{a_0 \sim \pi(s_0)} [\mathbb{E}_{s_1 \sim p(s_1|a_0, s_0)} [V^\pi(s_1)]] \\
&= \mathbb{E}_{a \sim \pi(s)} [r(s_0, a_0) + \gamma \mathbb{E}_{s_1 \sim p(s_1|a_0, s_0)} [V^\pi(s_1)]]
\end{aligned} \tag{1.2}$$

before we put s_1 to the right as the condition, it is stochastic, inside the $E_{s_1 \sim p(s_1|s_0, a_0)}$ scope it is deterministic, then we can get $V^\pi(s_1)$, as it needs the state to be deterministic.

The discrete formular is (get rid of the notation of time) so called **Bellman Equation**:

$$V^\pi(s) = \sum_{a \in \mathcal{A}} \pi(a \mid s) \left[r(s, a) + \gamma \sum_{s'} p(s' \mid s, a) V^\pi(s') \right], \forall s \in S \tag{1.3}$$

And if we write $r(s, a)$ as $\sum_r p(r \mid s, a) r$, then

$$p(r \mid s, a) = \sum_{s' \in \mathcal{S}} p(s', r \mid s, a)$$

We can also get

$$p(s' \mid s, a) = \sum_{r \in \mathcal{R}} p(s', r \mid s, a)$$

combined we get

$$V^\pi(s) = \sum_{a \in \mathcal{A}} \pi(a \mid s) \sum_{s' \in \mathcal{S}} \sum_{r \in \mathcal{R}} p(s', r \mid s, a) [r + \gamma V^\pi(s')] \tag{1.4}$$

If the reward depend solely on the next state s' , then

$$V^\pi(s) = \sum_{a \in \mathcal{A}} \pi(a \mid s) \sum_{s' \in \mathcal{S}} p(s' \mid s, a) [r(s') + \gamma V^\pi(s')] \tag{1.5}$$

Let

$$\begin{aligned}
r^\pi(s) &= \sum_{a \in \mathcal{A}} \pi(a \mid s) \sum_r p(r \mid s, a) r \\
p^\pi(s' \mid s) &= \sum_{a \in \mathcal{A}} \pi(a \mid s) p(s' \mid s, a)
\end{aligned} \tag{1.6}$$

rewirte 1.3 into the vector form:

$$V^\pi = r^\pi + \gamma P^\pi V^\pi \tag{1.7}$$

where $V^\pi = [V^\pi(s_1), \dots, V^\pi(s_n)]^\top \in \mathbb{R}^n$, $r^\pi = [r^\pi(s_1), \dots, r^\pi(s_n)]^\top \in \mathbb{R}^n$, and $P^\pi \in \mathbb{R}^{n \times n}$ with $P_{ij}^\pi = p^\pi(s_j \mid s_i)$.

1.4 Solving Value Function

Next, we need to solve the value function, first way is closed-form solution:

$$V^\pi = (I - \gamma P^\pi)^{-1} r^\pi$$

Some properties: $I - \gamma P^\pi$ is invertible, $(I - \gamma P^\pi)^{-1} \geq I$ which means every element of this inverse is nonnegative. For every vector $r \geq 0$, it holds that $(I - \gamma P^\pi)^{-1} r^\pi \geq r \geq 0$, so if $r_1 \geq r_2$, $(I - \gamma P^\pi)^{-1} r_1^\pi \geq (I - \gamma P^\pi)^{-1} r_2^\pi$

However, this method need to calculate the inverse of the matrix, that need some numerical algorithms. We can use a iterative solution:

$$V_{k+1} = r^\pi + \gamma P^\pi V_k$$

as $k \rightarrow \infty$, $V_k \rightarrow V^\pi = (I - \gamma P^\pi)^{-1} r^\pi$.

Proof Define the error as $\delta_k = V_k - V^\pi$, substitute $V_{k+1} = \delta_{k+1} + V^\pi$ and $V_k = \delta_k + V^\pi$ into the equation:

$$\delta_{k+1} + V^\pi = r^\pi + \gamma P^\pi (\delta_k + V^\pi)$$

Rearrange it:

$$\begin{aligned}
\delta_{k+1} &= r^\pi + \gamma P^\pi V^\pi + \gamma P^\pi \delta_k - V^\pi \\
&= \gamma P^\pi V^\pi + r^\pi + \gamma P^\pi \delta_k - V^\pi \\
&= \gamma P^\pi \delta_k
\end{aligned}$$

As a result, $\delta_{k+1} = \gamma P^\pi \delta_k = (\gamma P^\pi)^2 \delta_{k-1} = \dots = (\gamma P^\pi)^{k+1} \delta_0$. Since every entry of P^π is nonnegative and no greater than 1, and $\gamma < 1$, we have $\|(\gamma P^\pi)^{k+1}\| \rightarrow 0$ as $k \rightarrow \infty$, and the error $\|\delta_k\| \rightarrow 0$ as $k \rightarrow \infty$.

1.5 Action Value Function

Similarly to value function, **Action Value Function** is the value of an action at state s , from that state, take that action, the expected sum reward (return). We use $V^\pi(s)$ to denote value function, and $Q^\pi(s, a)$ to denote action value, their connection is:

$$V^\pi(s) = \sum_{a \in \mathcal{A}} \pi(a | s) Q^\pi(s, a) \quad (1.8)$$

The action value function is given as:

$$\begin{aligned} Q^\pi(s_0, a_0) &= r(s_0, a_0) + \mathbb{E}_{a_t \sim \pi(s_t)} [\sum_{t=1}^{\infty} \gamma^t r(s_t, a_t)] \\ &= r(s_0, a_0) + \gamma \mathbb{E}_{a_t \sim \pi(s_t)} [\sum_{t=1}^{\infty} \gamma^{t-1} r(s_t, a_t)] \\ &= r(s_0, a_0) + \gamma \mathbb{E}_{a_t \sim \pi(s_t)} [R(\tau')] \\ &= r(s_0, a_0) + \gamma \mathbb{E}_{s_1 \sim p(s_1 | s_0, a_0)} [\mathbb{E}_{a_t \sim \pi(s_t)} [R(\tau') | s_1]] \\ &= r(s_0, a_0) + \gamma \mathbb{E}_{s_1 \sim p(s_1 | s_0, a_0)} [V^\pi(s_1)] \\ &= r(s_0, a_0) + \gamma \mathbb{E}_{s_1 \sim p(s_1 | s_0, a_0)} [\sum_{a_1 \in \mathcal{A}} \pi(a_1 | s_1) Q^\pi(s_1, a_1)] \end{aligned} \quad (1.9)$$

Then the bellman equation of action value is:

$$\begin{aligned} Q^\pi(s, a) &= r(s, a) + \gamma \sum_{s'} p(s' | s, a) V^\pi(s') \\ &= r(s, a) + \gamma \sum_{s'} p(s' | s, a) \sum_{a' \in \mathcal{A}} \pi(a' | s') Q^\pi(s', a') \end{aligned} \quad (1.10)$$

Note that we can always write $r(s, a)$ as $\sum_r p(r | s, a) r$ if it is stochastic, and it follows the same notation in the book *Math of Reinforcement Learning*.

Rewrite 1.10 into vector form:

$$Q^\pi = \tilde{r} + \gamma P \Pi Q^\pi \quad (1.11)$$

where $\tilde{r}_{(s,a)} = \sum_r p(r | s, a) r$, $P_{(s,a),s'} = p(s' | s, a)$, $\Pi_{s', (s', a')} = \pi(a' | s')$.

1.6 Bellman Optimality Equation

Definition 1.1 (Optimal Policy)

If $V^{\pi_1}(s) \geq V^{\pi_2}(s), \forall s \in \mathcal{S}$, then π_1 is better than π_2 , if π_1 is better than all other policies, it is called **Optimal Policy** π^* .



Bellman Optimality Equation (BOE) is given by:

$$\begin{aligned} V(s) &= \max_{\pi(s) \in \Pi(s)} \sum_{a \in \mathcal{A}} \pi(a | s) (\sum_r p(r | s, a) r + \gamma \sum_{s'} p(s' | s, a) V(s')) \\ &= \max_{\pi(s) \in \Pi(s)} \sum_{a \in \mathcal{A}} \pi(a | s) Q(s, a) \end{aligned} \quad (1.12)$$

There are two unknowns in the equation, $V(s)$ and $\pi(a | s)$, we can first consider the right hand side, to compute the $\pi(a | s)$.

Example 1.1 Consider $\sum_1^3 c_i q_i$, where $c_1 + c_2 + c_3 = 1$ and they are all greater than 0, without loss of generality, we can assume $q_3 \geq q_1, q_2$, then the maximum is achieved when $c_3 = 1, c_1 = 0, c_2 = 0$. This is because:

$$q_3 = (c_1 + c_2 + c_3) q_3 = c_1 q_3 + c_2 q_3 + c_3 q_3 \geq c_1 q_1 + c_2 q_2 + c_3 q_3$$

Inspired by the example, since $\sum_a \pi(a | s) = 1$, we have:

$$\sum_{a \in \mathcal{A}} \pi(a | s) Q(s, a) \leq \sum_{a \in \mathcal{A}} \pi(a | s) \max_{a \in \mathcal{A}} Q(s, a) = \max_{a \in \mathcal{A}} Q(s, a)$$

where the equality is achieved when

$$\pi(a | s) = \begin{cases} 1, & a = a^*, \\ 0, & a \neq a^*. \end{cases}$$

here $a^* = \arg \max_{a \in \mathcal{A}} Q(s, a)$.

Then the matrix form of BOE is:

$$V = \max_{\pi \in \Pi} (r^\pi + \gamma P^\pi V) = f(V)$$

the r^π and P^π are the same before in normal Bellman equation 1.6.

In order to solve this nonlinear equation, we first need to introduce **Contraction Mapping** theorem or Fixed Point theorem:

Definition 1.2 (Contraction Mapping)

Consider function $f(x)$, where $x \in \mathbb{R}^d$ and $f : \mathbb{R}^d \rightarrow \mathbb{R}^d$. A point x^* is called a fixed point if $f(x^*) = x^*$, and the function is a contraction mapping if there exists $\gamma \in (0, 1)$ such that:

$$\|f(x_1) - f(x_2)\| \leq \gamma \|x_1 - x_2\|, \forall x_1, x_2 \in \mathbb{R}^d$$



The relation between a fixed point and the contraction property is characterized by:

Theorem 1.1 (Banach's Fixed Point Theorem)

For any equation that has the form $x = f(x)$ where x and $f(x)$ are real vectors, if f is a contraction mapping, then:

1. *Existence:* There exists a fixed point x^* such that $f(x^*) = x^*$.
2. *Uniqueness:* There exists a unique fixed point x^* such that $f(x^*) = x^*$.
3. *Algorithm:* For any initial point x_0 , the sequence $x_{k+1} = f(x_k)$ converges to the fixed point x^* . Moreover, the convergence rate is exponentially fast.



The proof of the theorem can be found in the book, it is based on Cauchy sequence. Then we need to show the right hand side of the BOE is a contraction mapping:

Theorem 1.2 (Contraction Property of right-hand side of BOE)

For any $V_1, V_2 \in \mathbb{R}^{|S|}$, it holds that:

$$\|f(V_1) - f(V_2)\|_\infty \leq \gamma \|V_1 - V_2\|_\infty$$

where $\gamma \in (0, 1)$ is the discount factor, $\|\cdot\|_\infty$ is the maximum norm, which is the maximum absolute value of the elements of a vector.



Proof

$$f(V_1) = \max_{\pi} (r^\pi + \gamma P^\pi V_1) = r^{\pi_1^*} + \gamma P^{\pi_1^*} V_1 \geq r^{\pi_2^*} + \gamma P^{\pi_2^*} V_1$$

$$f(V_2) = \max_{\pi} (r^\pi + \gamma P^\pi V_2) = r^{\pi_2^*} + \gamma P^{\pi_2^*} V_2 \geq r^{\pi_1^*} + \gamma P^{\pi_1^*} V_2$$

where \geq is elementwise comparison, as a result:

$$\begin{aligned} f(V_1) - f(V_2) &= (r^{\pi_1^*} - r^{\pi_2^*}) + \gamma P^{\pi_1^*} V_1 - \gamma P^{\pi_2^*} V_2 \\ &\leq (r^{\pi_1^*} - r^{\pi_1^*}) + \gamma P^{\pi_1^*} V_1 - \gamma P^{\pi_1^*} V_2 \\ &= \gamma P^{\pi_1^*} (V_1 - V_2) \end{aligned}$$

similarly we can get $f(V_2) - f(V_1) \leq \gamma P^{\pi_2^*} (V_2 - V_1)$, so we have:

$$\gamma P^{\pi_2^*} (V_1 - V_2) \leq f(V_1) - f(V_2) \leq \gamma P^{\pi_1^*} (V_1 - V_2)$$

define

$$z = \max \{ |\gamma P^{\pi_1^*} (V_1 - V_2)|, |\gamma P^{\pi_2^*} (V_1 - V_2)| \} \in \mathbb{R}^{|S|}$$

all the operations are elementwise, $z \geq 0$, then we have:

$$-z \leq \gamma P^{\pi_2^*} (V_1 - V_2) \leq f(V_1) - f(V_2) \leq \gamma P^{\pi_1^*} (V_1 - V_2) \leq z$$

which implies:

$$|f(V_1) - f(V_2)| \leq z$$

it then follows that:

$$\|f(V_1) - f(V_2)\|_\infty \leq \|z\|_\infty \quad (1.13)$$

suppose z_i, p_i^T, q_i^T are i th entry of $z, P^{\pi_1^*}, P^{\pi_2^*}$, then:

$$z_i = \max \{ |\gamma p_i^T (V_1 - V_2)|, |\gamma q_i^T (V_1 - V_2)| \}$$

since p_i sums up to 1 and nonnegative, we have:

$$|p_i^T (V_1 - V_2)| \leq p_i^T |V_1 - V_2| \leq \|V_1 - V_2\|_\infty$$

similarly we have $|q_i^T (V_1 - V_2)| \leq \|V_1 - V_2\|_\infty$, therefore $z_i \leq \gamma \|V_1 - V_2\|_\infty$, and hence

$$\|z\|_\infty = \max_i |z_i| \leq \gamma \|V_1 - V_2\|_\infty$$

Substitute it back to 1.13, we have:

$$\|f(V_1) - f(V_2)\|_\infty \leq \gamma \|V_1 - V_2\|_\infty$$

Then we can use this to solve an optimal policy from the BOE. Since $V^* = \max_{\pi \in \Pi} (r^\pi + \gamma P^\pi V^*)$, so it is clearly a fixed point.

Theorem 1.3 (Existence, Uniqueness and Algorithm of Optimal Policy)

The optimal policy V^* exists and is unique, and the sequence $V_{k+1} = f(V_k)$ converges to the optimal policy V^* exponentially fast given any initial guess V_0 with the iteration algorithm:

$$V_{k+1} = f(V_k) = \max_{\pi \in \Pi} (r^\pi + \gamma P^\pi V_k)$$



The proof follows the proof of the contraction mapping theorem, and the iteration algorithm is called **Value Iteration**. Once we have the optimal value function V^* , we can get the optimal policy π^* by:

$$\pi^* = \arg \max_{\pi \in \Pi} (r^\pi + \gamma P^\pi V^*) \quad (1.14)$$

substitute this into the BOE yields $V^* = r^{\pi^*} + \gamma P^{\pi^*} V^*$, which is the optimal value function.

Theorem 1.4 (Optimality of Value Function and Policy)

The solution V^* and the policy π^* are optimal, i.e., for any other policy $\pi \in \Pi$, it holds that:

$$V^* = V^{\pi^*} \geq V^\pi$$



Proof

$$V^* - V^\pi = r^{\pi^*} + \gamma P^{\pi^*} V^* - r^\pi - \gamma P^\pi V^\pi \geq r^{\pi^*} + \gamma P^{\pi^*} V^* - r^\pi - \gamma P^{\pi^*} V^\pi = \gamma P^{\pi^*} (V^* - V^\pi)$$

iteratively we have $V^* - V^\pi \geq \lim_{n \rightarrow \infty} (\gamma P^{\pi^*})^n (V^* - V^\pi)$.

Theorem 1.5 (Greedy optimal policy)

For any $s \in \mathcal{S}$, the deterministic greedy policy:

$$\pi(a | s) = \begin{cases} 1, & a = a^*, \\ 0, & a \neq a^*. \end{cases} \quad (1.15)$$

is an optimal policy, where $a^* = \arg \max_{a \in \mathcal{A}} Q^*(s, a)$, where

$$Q^*(s, a) = \sum_{r \in \mathcal{R}} p(r | s, a) r + \gamma \sum_{s'} p(s' | s, a) V^*(s')$$

remember that even though V^* is unique, the optimal policy may not be unique, and there always exists a greedy optimal policy.




We can talk about the impact of the reward values:

Theorem 1.6 (Optimal policy invariance)

If every reward $r \in \mathcal{R}$ is changed by an affine transformation to $\alpha r + \beta$, where $\alpha, \beta \in \mathbb{R}$ and $\alpha > 0$, then the corresponding optimal state value V' is also an affine transformation of V^* :

$$V' = \alpha V^* + \frac{\beta}{1 - \gamma} \mathbf{1}$$

Consequently, the optimal policy derived from V' is invariant to the affine transformation of the reward values. 

And with the discount factor γ , the optimal policy will not take any meaningless detour.

1.7 Value Iteration and Policy Iteration

The algorithm of *Value Iteration* is:

$$V_{k+1} = \max_{\pi \in \Pi} (r^\pi + \gamma P^\pi V_k)$$

and there are two steps in one iteration, first one is called *Policy Update*:

$$\pi_{k+1} = \arg \max_{\pi \in \Pi} (r^\pi + \gamma P^\pi V_k)$$

the second one is called *Value Update*:

$$V_{k+1} = r^{\pi_{k+1}} + \gamma P^{\pi_{k+1}} V_k$$

With the elementwise form, the policy update is (if there are same actions that takes the maximum, we can choose any of them):

$$\pi_{k+1}(s) = \arg \max_{\pi(s) \in \Pi(s)} \sum_{a \in \mathcal{A}} \pi(a | s) \left(\sum_r p(r | s, a) r + \gamma \sum_{s'} p(s' | s, a) V_k(s') \right)$$

Then the value update is:

$$V_{k+1}(s) = \sum_{a \in \mathcal{A}} \pi_{k+1}(a | s) \left(\sum_r p(r | s, a) r + \gamma \sum_{s'} p(s' | s, a) V_k(s') \right)$$

We can use the greedy deterministic policy, which is then:

$$V_{k+1}(s) = \max_a Q_k(s, a)$$

We should know that V_k is not a state value though it converges to the optimal state value, it is not ensured to satisfy the Bellman equation. So the Q_k is also not a action value, they are all intermediate values.

The algorithm of *Policy Iteration* has also two steps, first one is called *Policy Evaluation*:

$$V^{\pi_k} = r^{\pi_k} + \gamma P^{\pi_k} V^{\pi_k}$$


second one is called *Policy Improvement*:

$$\pi_{k+1} = \arg \max_{\pi \in \Pi} (r^{\pi_k} + \gamma P^{\pi_k} V^{\pi_k})$$

Here comes the first question, how to solve the policy evaluation? We can use the iterative method introduced in 1.4, and this results an iteration algorithm inside an iteration algorithm. We will not do infinite iteration here, so the V^{π_k} will not be the exact solution, would this cause problem? No, see the truncated policy iteration algorithm.

And the second question, why π_{k+1} is better than π_k ?

Lemma 1.1 (Policy Improvement)

If $\pi_{k+1} = \arg \max_{\pi \in \Pi} (r^{\pi_k} + \gamma P^{\pi_k} V^{\pi_k})$, then $V^{\pi_{k+1}}(s) \geq V^{\pi_k}(s), \forall s \in \mathcal{S}$. 

Proof We know that:

$$r^{\pi_{k+1}} + \gamma P^{\pi_{k+1}} V^{\pi_{k+1}} \geq r^{\pi_k} + \gamma P^{\pi_k} V^{\pi_k}$$

Then

$$\begin{aligned}
 V^{\pi_k} - V^{\pi_{k+1}} &= r^{\pi_k} + \gamma P^{\pi_k} V^{\pi_k} - (r^{\pi_{k+1}} + \gamma P^{\pi_{k+1}} V^{\pi_{k+1}}) \\
 &\leq r^{\pi_{k+1}} + \gamma P^{\pi_{k+1}} V^{\pi_k} - (r^{\pi_{k+1}} + \gamma P^{\pi_{k+1}} V^{\pi_{k+1}}) \\
 &\leq \gamma P^{\pi_{k+1}} (V^{\pi_k} - V^{\pi_{k+1}})
 \end{aligned}$$

again iteratively we have:

$$V^{\pi_k} - V^{\pi_{k+1}} \leq \lim_{n \rightarrow \infty} (\gamma P^{\pi_{k+1}})^n (V^{\pi_k} - V^{\pi_{k+1}}) = 0$$

Theorem 1.7 (Convergence of policy iteration)

The state value sequence $\{V^{\pi_k}\}_{k=0}^{\infty}$ converges to the optimal state value V^* , and the policy sequence $\{\pi_k\}_{k=0}^{\infty}$ converges to the optimal policy π^* in policy iteration algorithm.

Proof We introduce another sequence $\{V_k\}_{k=0}^{\infty}$ generated by:

$$V_{k+1} = f(V_k) = \max_{\pi} (r^{\pi} + \gamma P^{\pi} V_k)$$

which is exactly the value iteration algorithm, and we know that $V_k \rightarrow V^*$ as $k \rightarrow \infty$. For $k = 0$, we can always find a V_0 such that $V_0 \leq V^{\pi_0}$, then we use induction: for $k \geq 0$, if $V_k \leq V^{\pi_k}$, then for $k + 1$:

$$\begin{aligned}
 V^{\pi_{k+1}} - V_{k+1} &= (r^{\pi_{k+1}} + \gamma P^{\pi_{k+1}} V^{\pi_{k+1}}) - \max_{\pi} (r^{\pi} + \gamma P^{\pi} V_k) \\
 &\geq (r^{\pi_{k+1}} + \gamma P^{\pi_{k+1}} V_k) - \max_{\pi} (r^{\pi} + \gamma P^{\pi} V_k) \\
 &\quad (\text{because } V^{\pi_{k+1}} \geq V_k \text{ by Lemma 1.1 and } P^{\pi_{k+1}} \geq 0) \\
 &= (r^{\pi'_k} + \gamma P^{\pi'_k} V_k) - (r^{\pi'_k} + \gamma P^{\pi'_k} V_k) \\
 &\quad (\text{suppose } \pi'_k = \arg \max_{\pi} (r^{\pi} + \gamma P^{\pi} V_k)) \\
 &\geq (r^{\pi'_k} + \gamma P^{\pi'_k} V_k) - (r^{\pi'_k} + \gamma P^{\pi'_k} V_k) \\
 &\quad (\text{because } \pi_{k+1} = \arg \max_{\pi} (r^{\pi} + \gamma P^{\pi} V_k)) \\
 &= \gamma P^{\pi'_k} (V_{k+1} - V_k) \geq 0
 \end{aligned}$$

Since V_k converges to V^* , V^{π_k} is also converges to V^* .

Then the elementwise form of the policy iteration algorithm is, policy evaluation:

$$V_{(j+1)}^{\pi_k}(s) = \sum_{a \in \mathcal{A}} \pi_{k+1}(a | s) \left(\sum_r p(r | s, a) r + \gamma \sum_{s'} p(s' | s, a) V_{(j)}^{\pi_k}(s') \right)$$

and policy improvement:

$$\pi_{k+1}(s) = \arg \max_{\pi(s) \in \Pi(s)} \sum_{a \in \mathcal{A}} \pi(a | s) \left(\sum_r p(r | s, a) r + \gamma \sum_{s'} p(s' | s, a) V^{\pi_k}(s') \right)$$

Next we introduce **Truncated Policy Iteration** algorithm, which is a combination of value iteration and policy iteration. We will see that the value iteration and policy iteration algorithms are two special cases of the truncated policy iteration algorithm.

If we start from $V_0^{\pi_1} = V_0$, we get:

$$\begin{array}{lll}
 & & V_{\pi_1}^{(0)} = V_0 \\
 \text{value iteration} & \leftarrow V_1 \leftarrow & V_{(1)}^{\pi_1} = r^{\pi_1} + \gamma P^{\pi_1} V_{(0)}^{\pi_1} \\
 & & V_{(2)}^{\pi_1} = r^{\pi_1} + \gamma P^{\pi_1} V_{(1)}^{\pi_1} \\
 & & \vdots \\
 \text{truncated policy iteration} & \leftarrow \bar{V}_1 \leftarrow & V_{(j)}^{\pi_1} = r^{\pi_1} + \gamma P^{\pi_1} V_{(j-1)}^{\pi_1} \\
 & & \vdots \\
 \text{policy iteration} & \leftarrow V^{\pi_1} \leftarrow & V_{(\infty)}^{\pi_1} = r^{\pi_1} + \gamma P^{\pi_1} V_{(\infty)}^{\pi_1}
 \end{array}$$

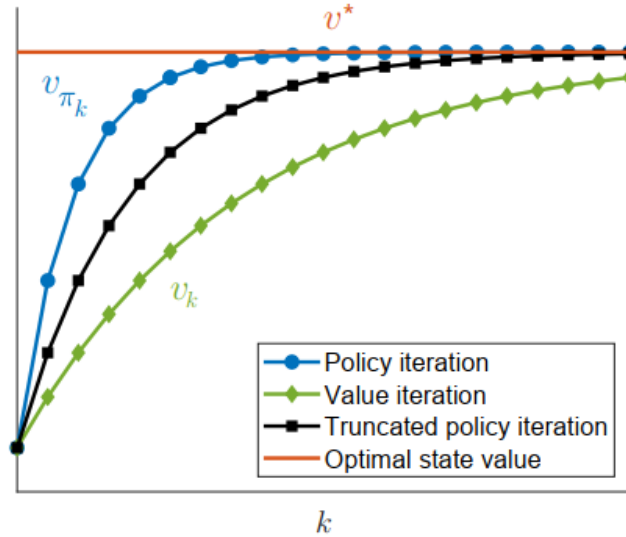


Figure 1.1: 3 iteration methods

Proposition 1.1 (Value Improvement)

In the policy evaluation step if the initial guess is selected as $V_{(0)}^{\pi_k} = V^{\pi_{k-1}}$, then it holds that:

$$V_{(j+1)}^{\pi_k} \geq V_{(j)}^{\pi_k}$$

Proof

$$V_{(j+1)}^{\pi_k} - V_{(j)}^{\pi_k} = \gamma P^{\pi_k} (V_{(j)}^{\pi_k} - V_{(j-1)}^{\pi_k}) = \dots = \gamma^j P^{\pi_k} (V_{(1)}^{\pi_k} - V_{(0)}^{\pi_k})$$

we have

$$V_{(1)}^{\pi_k} = r^{\pi_k} + \gamma P^{\pi_k} V^{\pi_{k-1}} \geq r^{\pi_{k-1}} + \gamma P^{\pi_{k-1}} V^{\pi_{k-1}} = V^{\pi_{k-1}} = V_{(0)}^{\pi_k}$$

where the inequality is due to $\pi_k = \arg \max_{\pi} (r^{\pi} + \gamma P^{\pi} V^{\pi_{k-1}})$, substitute $V_{(1)}^{\pi_k} \geq V_{(0)}^{\pi_k}$ into the first equation, we have $V_{(j+1)}^{\pi_k} \geq V_{(j)}^{\pi_k}$.

But we have to note that this is based on the assumption that the initial guess is $V_{(0)}^{\pi_k} = V^{\pi_{k-1}}$, and $V^{\pi_{k-1}}$ is not always available, like in the truncated policy iteration algorithm, we do not do the iteration until convergence, but only do a few steps, so the initial value is approximation. And in Deep Reinforcement Learning, we do not have the exact value function, even if we have it may not be transferable or meaningful.

1.8 Monte Carlo Method

If we sample from a independent and identically distributed (i.i.d.) distribution, we can use the Monte Carlo method to do mean estimate, as value function and action value function are all mean estimates.

In the policy iteration, the action value lies in the core of these two steps, we first use Bellman equation to compute the value function, then use it to compute action value, and then update the policy. If we do not have a model, then the definition of action value is:

$$Q^{\pi_k}(s, a) = \mathbb{E}[G_t | S_t = s, A_t = a] = \mathbb{E}[R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots | S_t = s, A_t = a]$$

we use R to denote the stochastic reward, and G_t is the return at time t .

If we run the policy for n episodes then we get:

$$Q^{\pi_k}(s, a) = \mathbb{E}[G_t | S_t = s, A_t = a] \approx \frac{1}{n} \sum_{i=1}^n G_{(i)}^{\pi_k}(s, a) \quad (1.16)$$

instead of calculation of value function, we directly get action value from samples, this is the simplest Monte Carlo method (we call it MC Basic).

In the example of the book, we can see that if the length of the episodes are too short, it will get zero action value, because it can not get to the target. This is caused by the *Sparse Reward*.

We next show how to efficiently use the samples. Every time a state-action pair appears in an episodes, it is called a *visit*. Given an episodes:

$$s_1 \xrightarrow{a_2} s_2 \xrightarrow{a_4} s_1 \xrightarrow{a_2} s_2 \xrightarrow{a_3} s_5 \xrightarrow{a_1} \dots$$

1. Initial visit, only the first pair is evaluated, like in the MC Basic, which is not efficient.
2. First-visit, only use the samples to estimate the action value of the first visit.
3. Every visit, use the samples to estimate the action value of every visit, which is more efficient, but may cause bias (because they are subsequence, there are correlation).

Then we consider how to efficiently update the policy.

1. In the policy evaluation step, to collect all the episodes starting from the same state-action pair and then approximate the action value using the average return of these episodes, which is in MC Basic. The drawback of this strategy is that the agent must wait until all the episodes have been collected before the estimate can be updated.
2. Use the return of a single episode to approximate the corresponding action value. In this way, we can immediately obtain a rough estimate when we receive an episode. Then, the policy can be improved in an episode-by-episode fashion. In fact, this strategy falls into the scope of *Generalized Policy Iteration* introduced in the last chapter, we can still update the policy even if the value estimate is not sufficiently accurate.

So the new algorithm is called *MC Exploring Starts*, based on the exploring starts condition, which means that the agent can start from any state and take any action (MC Basic assume this too).

Algorithm 1: MC Exploring Starts (efficient variant of MC Basic)

Input: Initial policy $\pi_0(a|s)$, $Q(s, a)$, $\text{Returns}(s, a) = 0$, $\text{Num}(s, a) = 0, \forall (s, a)$

```

1 foreach episode do
2   Episode generation: Select  $(s_0, a_0)$  (exploring-starts), then generate  $s_0, a_0, r_1, \dots, s_{T-1}, a_{T-1}, r_T$ 
   following  $\pi$ 
3    $g \leftarrow 0$ 
4   for  $t = T - 1, T - 2, \dots, 0$  do
5      $g \leftarrow \gamma g + r_{t+1}$ 
6      $\text{Returns}(s_t, a_t) \leftarrow \text{Returns}(s_t, a_t) + g$ 
7      $\text{Num}(s_t, a_t) \leftarrow \text{Num}(s_t, a_t) + 1$ 
8      $Q(s_t, a_t) \leftarrow \frac{\text{Returns}(s_t, a_t)}{\text{Num}(s_t, a_t)}$ 
9      $\pi(a|s_t) \leftarrow \begin{cases} 1 & \text{if } a = \arg \max_a Q(s_t, a), \\ 0 & \text{otherwise} \end{cases}$ 
10  end
11 end
```

We can see that the algorithm starts from the last state of the episode, and then update the action value and policy in a **reverse order**.

However, the exploring starts condition is not always satisfied, for example in the real world involving physical interactions, we need to remove this condition. A policy is **soft** if it has a positive probability of taking any action at any state, with a soft policy, a single episode that is sufficiently long can visit every state-action pair many times. Thus we do not need the exploring starts condition.

One type of soft policy is *ϵ -greedy Policy*:

$$\pi(a | s) = \begin{cases} 1 - \frac{\epsilon}{|\mathcal{A}(s)|} (|\mathcal{A}(s)| - 1), & \text{for the greedy action,} \\ \frac{\epsilon}{|\mathcal{A}(s)|}, & \text{for the other actions.} \end{cases} \quad (1.17)$$

Where $\pi \in \Pi_\epsilon$ is the set of all ϵ -greedy policies, and $|\mathcal{A}(s)|$ is the number of actions available at state s . The probability of taking the greedy action is always greater than that of taking any other action.

With this we can change the policy improvement step:

$$\pi(a | s) = \begin{cases} 1 - \frac{\epsilon}{|\mathcal{A}(s)|} (|\mathcal{A}(s)| - 1), & a = a^*, \\ \frac{\epsilon}{|\mathcal{A}(s)|}, & a \neq a^*. \end{cases} \quad (1.18)$$

So the algorithm is:

Algorithm 2: MC ϵ -greedy (a variant of MC Exploring Starts)

Input: Initial policy $\pi_0(a|s)$, $Q(s, a)$, $\text{Returns}(s, a) = 0$, $\text{Num}(s, a) = 0, \forall (s, a). \epsilon \in (0, 1]$

```

1 foreach episode do
2   Episode generation: Select  $(s_0, a_0)$ , then generate  $s_0, a_0, r_1, \dots, s_{T-1}, a_{T-1}, r_T$  following  $\pi$ 
3    $g \leftarrow 0$ 
4   for  $t = T - 1, T - 2, \dots, 0$  do
5      $g \leftarrow \gamma g + r_{t+1}$ 
6      $\text{Returns}(s_t, a_t) \leftarrow \text{Returns}(s_t, a_t) + g$ 
7      $\text{Num}(s_t, a_t) \leftarrow \text{Num}(s_t, a_t) + 1$ 
8      $Q(s_t, a_t) \leftarrow \frac{\text{Returns}(s_t, a_t)}{\text{Num}(s_t, a_t)}$ 
9      $\pi(a|s_t) \leftarrow \begin{cases} 1 - \frac{\epsilon}{|\mathcal{A}(s)|} (|\mathcal{A}(s)| - 1) & a = a^*, \\ \frac{\epsilon}{|\mathcal{A}(s)|} & a \neq a^*. \end{cases}$ 
10  end
11 end
```

The policy is optimal in Π_ϵ , not in Π , so it is not guaranteed to be optimal, but it is still a good policy (if ϵ is sufficiently small).

1.9 Stochastic Approximation

For now, we only learned non-incremental algorithms, which may take a logn time to compute, not good for online learning. Suppose $w_k = \frac{1}{k-1} \sum_{i=1}^{k-1} x_i, k = 2, 3, \dots$, then:

$$w_{k+1} = \frac{1}{k} \sum_{i=1}^k x_i = \frac{k-1}{k} w_k + \frac{1}{k} x_k$$

this update the mean estimate incrementally. We use

$$w_{k+1} = w_k - \alpha_k (w_k - x_k) \quad (1.19)$$

α_k may not be given, though above we give it as $\frac{1}{k}$.

Given a noisy observation of $g(w)$:

$$\tilde{g}(w, \eta) = g(w) + \eta$$

where η is a noise, we can use the stochastic approximation method **Robbins-Monro** algorithm to update the estimate of w :

$$w_{k+1} = w_k - \alpha_k \tilde{g}(w_k, \eta_k), \alpha_k > 0 \quad (1.20)$$

Why this algorithm converge?

Theorem 1.8 (Robbins-Monro theorem)

If

1. $0 \leq c_1 \leq \nabla_w g(w) \leq c_2 < \infty, \forall w$; So g is monotonic and bounded;
2. $\sum_{k=1}^{\infty} \alpha_k = \infty$ and $\sum_{k=1}^{\infty} \alpha_k^2 < \infty$; It should converge to zero, but not too fast, one common choice is $\alpha_k = \frac{1}{k}$;
3. $\mathbb{E}[\eta_k | \mathcal{H}_k] = 0$ and $\mathbb{E}[\eta_k^2 | \mathcal{H}_k] < \infty$;

where $\mathcal{H}_k = \{w_k, w_{k-1}, \dots\}$, then w_k converges to the root of $g(w)$ almost surely, i.e., $\lim_{k \rightarrow \infty} w_k = w^*$, where $g(w^*) = 0$.

Then let $g(w) = w - \mathbb{E}[X]$, $\tilde{g}(w, \eta) = w - x + \eta$, $\eta = \mathbb{E}[X] - x$, we can use the Robbins-Monro algorithm to prove that with careful choice of α_k , the estimate w_k converges to the mean $\mathbb{E}[X]$ almost surely.

Theorem 1.9 (Dvoretzky's Theorem)

Consider a stochastic sequence:

$$\Delta_{k+1} = (1 - \alpha_k)\Delta_k + \beta_k \eta_k$$

where $\alpha_k, \beta_k \geq 0$, then Δ_k converges to 0 almost surely, if:

1. $\sum_{k=1}^{\infty} \alpha_k = \infty$, $\sum_{k=1}^{\infty} \alpha_k^2 < \infty$, and $\sum_{k=1}^{\infty} \beta_k^2 < \infty$ uniformly almost surely;
2. $\mathbb{E}[\eta_k | \mathcal{H}_k] = 0$ and $\mathbb{E}[\eta_k^2 | \mathcal{H}_k] < \infty$ almost surely.

where $\mathcal{H}_k = \{\Delta_k, \Delta_{k-1}, \dots, \eta_{k-1}, \dots, \alpha_{k-1}, \dots, \beta_{k-1}, \dots\}$.

In RM, α_k, β_k are deterministic, but in Dvoretzky's theorem, they can be stochastic, and is useful in cases they are function of Δ_k .

Proof Please refer to the book, page 110.

And we can also use this theorem to prove the convergence of the mean estimation, Robbins-Monro algorithm as well, and please refer to the book too.

Theorem 1.10 (Exentsion of Dvoretzky's Theorem)

Consider a finite set \mathcal{S} of real numbers, a stochastic process:

$$\Delta_{k+1}(s) = (1 - \alpha_k(s))\Delta_k(s) + \beta_k(s)\eta_k(s)$$

$\Delta_k(s)$ converges to 0 almost surely, if:

1. $\sum_{k=1}^{\infty} \alpha_k(s) = \infty$, $\sum_{k=1}^{\infty} \alpha_k^2(s) < \infty$, $\sum_{k=1}^{\infty} \beta_k^2(s) < \infty$ and $\mathbb{E}[\beta_k(s) | \mathcal{H}_k] \leq \mathbb{E}[\alpha_k(s) | \mathcal{H}_k]$ uniformly almost surely.
2. $\|\mathbb{E}[\eta_k | \mathcal{H}_k]\|_{\infty} \leq \gamma \|\Delta_k\|_{\infty}$, where $\gamma \in (0, 1)$
3. $\text{var}[\eta_k | \mathcal{H}_k] \leq C(1 + \|\Delta_k\|_{\infty})^2$, where C is a constant.

where $\mathcal{H}_k = \{\Delta_k, \Delta_{k-1}, \dots, \eta_{k-1}, \dots, \alpha_{k-1}, \dots, \beta_{k-1}, \dots\}$.

- The variable s can be viewed as an index, in the context of reinforcement learning, it indicates a state or a state-action pair. The norm is maximum norm.
- This theorem only requires that the expectation and variance are bounded by the error. The convergence requires the conditions are valid for every s .

We then show that **stochastic gradient descent** (SGD) is a special case of the Robbins-Monro algorithm, and mean estimation is a special case of the SGD. The preliminary knowledge of gradient descent, convex function and the convergence of the GD please refer to the book appendix D.

Given a optimization problem:

$$\min_w J(w) = \mathbb{E}[f(w, X)]$$

where f is a function of parameter w and random vairable X , we can use the gradient descent to solve it:

$$w_{k+1} = w_k - \alpha_k \nabla_{w_k} J(w_k) = w_k - \alpha_k \mathbb{E}[\nabla_w f(w_k, X_k)] \quad (1.21)$$

The expectation is not computable, so we collect a larger number of data:

$$w_{k+1} = w_k - \frac{\alpha_k}{n} \sum_{i=1}^n \nabla_w f(w_k, x_i)$$

And we if only use one sample at a time, we can use the stochastic gradient descent (SGD):

$$w_{k+1} = w_k - \alpha_k \nabla_w f(w_k, x_k) \quad (1.22)$$

An interesting **convergence pattern** is that it behaves similarly to the regular gradient descent algorithm when the estimate w_k is far from the optimal solution w^* . Only when they are close, does the convergence of SGD exhibit more randomness.

$$\begin{aligned}
 \delta_k &= \frac{|\nabla_w f(w_k, x_k) - \mathbb{E}[\nabla_w f(w_k, X)]|}{|\mathbb{E}[\nabla_w f(w_k, X)]|} \\
 &= \frac{|\nabla_w f(w_k, x_k) - \mathbb{E}[\nabla_w f(w_k, X)]|}{|\mathbb{E}[\nabla_w f(w_k, X)] - \mathbb{E}[\nabla_w f(w^*, X)]|} \quad \text{as } \mathbb{E}[\nabla_w f(w^*, X)] = 0 \\
 &= \frac{|\nabla_w f(w_k, x_k) - \mathbb{E}[\nabla_w f(w_k, X)]|}{|\mathbb{E}[\nabla_w^2 f(w_k, X)](w_k - w^*)|} \quad (\text{by mean value theorem}) \\
 &\leq \frac{|\nabla_w f(w_k, x_k) - \mathbb{E}[\nabla_w f(w_k, X)]|}{c |w_k - w^*|}
 \end{aligned}$$

Above we use stochastic formulation, but we can also use the deterministic formulation, the equations are the same, but samples are not random, they are fixed $\{x_i\}_{i=1}^n$. So does the same equations are SGD too? The answer is yes, because if we set X be a random variable that $P(X = x_i) = \frac{1}{n}$, then $\min_w J(w) = \frac{1}{n} \sum f(w, x_i) = \mathbb{E}[f(w, X)]$. It means if we uniformly and independently sample x_k from the $\{x_i\}_{i=1}^n$, it is still SGD.

Theorem 1.11 (Convergence of SGD)

w_k converges to the optimal solution w^* almost surely, where $\mathbb{E}[\nabla_w f(w^*, X)] = 0$, if:

1. $0 < c_1 \leq \nabla_w^2 f(w, X) \leq c_2$;
2. $\sum_{k=1}^{\infty} \alpha_k = \infty$ and $\sum_{k=1}^{\infty} \alpha_k^2 < \infty$;
3. $\{x_k\}_{k=1}^{\infty}$ is i.i.d.



Proof Let

$$g(w) = \nabla_w J(w) = \mathbb{E}[\nabla_w f(w, X)]$$

Then SGD aims to find the root of $g(w) = 0$, and $\tilde{g} = \nabla_w f(w, x)$ where x is a sample of X :

$$\tilde{g}(w, \eta) = \mathbb{E}[\nabla_w f(w, X)] + \nabla_w f(w, x) - \mathbb{E}[\nabla_w f(w, X)] = \mathbb{E}[\nabla_w f(w, X)] + \eta$$

Then the RM is:

$$w_{k+1} = w_k - \alpha_k \tilde{g}(w, \eta) = w_k - \alpha_k \nabla_w f(w, x)$$

Then we show it follows 3 conditions in 1.8:

1. $\nabla_w g(w) = \nabla_w \mathbb{E}[\nabla_w f(w, X)] = \mathbb{E}[\nabla_w^2 f(w, X)]$, then because $0 < c_1 \leq \nabla_w^2 f(w, X) \leq c_2$, $c_1 \leq \nabla_w g(w) \leq c_2$.
2. second conditions are the same
3. since $\{x_k\}_{k=1}^{\infty}$ is i.i.d, $\mathbb{E}_{x_k}[\nabla_w f(w, x_k)] = \mathbb{E}[\nabla_w f(w, X)]$ for all k , therefore

$$\mathbb{E}[\eta_k | \mathcal{H}_k] = \mathbb{E}[\nabla_w f(w_k, x_k) - \mathbb{E}[\nabla_w f(w_k, X)] | \mathcal{H}_k] = 0$$

as $\mathbb{E}[\nabla_w f(w_k, x_k) | \mathcal{H}_k] = \mathbb{E}_{x_k}[\nabla_w f(w_k, x_k)]$ because of independency, $\mathbb{E}[\mathbb{E}[\nabla_w f(w_k, X)] | \mathcal{H}_k] = \mathbb{E}[\nabla_w f(w_k, X)]$ because $\mathbb{E}[\nabla_w f(w_k, X)]$ is function of w_k . Same we can show $\mathbb{E}[\eta_k^2 | \mathcal{H}_k] < \infty$.

1.10 Temporal-Difference Methods

The algorithm is:

$$v_{t+1}(s_t) = v_t(s_t) - \alpha_t(s_t)[v_t(s_t) - (r_{t+1} + \gamma v_t(s_{t+1}))] \quad (1.23)$$

and $v_{t+1}(s) = v_t(s)$ for all $s \neq s_t$, where $t = 0, 1, 2, \dots$. Here $v_t(s_t)$ is the estimate of $V^\pi(s_t)$, $\alpha_t(s_t)$ is the learning rate. Only value that are visited will update!

$$V^\pi(s) = \mathbb{E}[R_{t+1} + \gamma G_{t+1} | S_t = s] = \mathbb{E}[R_{t+1} + \gamma V^\pi(S_{t+1}) | S_t = s]$$

Then let

$$g(V^\pi(s_t)) = V^\pi(s_t) - \mathbb{E}[R_{t+1} + \gamma V^\pi(S_{t+1}) \mid S_t = s] = 0$$

But we will only get samples so

$$\begin{aligned} \tilde{g}(V^\pi(s_t)) &= V^\pi(s_t) - [r_{t+1} + \gamma V^\pi(s_{t+1})] \\ &= \underbrace{(V^\pi(s_t) - \mathbb{E}[R_{t+1} + \gamma V^\pi(S_{t+1}) \mid S_t = s_t])}_{g(V^\pi(s_t))} \\ &\quad + \underbrace{(\mathbb{E}[R_{t+1} + \gamma V^\pi(S_{t+1}) \mid S_t = s_t] - [r_{t+1} + \gamma V^\pi(s_{t+1})])}_{\eta}. \end{aligned}$$

then we get RM algorithm as:

$$v_{t+1}(s_t) = v_t(s_t) - \alpha_t(s_t)[v_t(s_t) - (r_{t+1} + \gamma V^\pi(s_{t+1}))]$$

the only difference is V^π replaces $v_t(s_{t+1})$. And we will show this replacement does not affect the convergence of TD in 1.12.

$$\underbrace{v_{t+1}(s_t)}_{\text{new estimate}} = \underbrace{v_t(s_t)}_{\text{current estimate}} - \alpha_t(s_t) \underbrace{[v_t(s_t) - (r_{t+1} + \gamma v_t(s_{t+1}))]}_{\text{TD error } \delta_t \text{ TD target } \bar{v}_t}$$

\bar{v}_t is called TD target because:

$$\begin{aligned} v_{t+1}(s_t) - \bar{v}_t &= [1 - \alpha_t(s_t)][v_t(s_t) - \bar{v}_t] \\ &\rightarrow \text{Get the absolute value} \\ |v_{t+1}(s_t) - \bar{v}_t| &= |1 - \alpha_t(s_t)||v_t(s_t) - \bar{v}_t| \\ &< |v_t(s_t) - \bar{v}_t| \end{aligned}$$

The expectation of TD error is 0 if $v_t(s_t) = V^\pi(s_t)$, so it reflects the discrepancy between the estimate and the true state value, it can be interpreted as *innovation*, which indicates new information obtained from the experience sample. A comparison of TD and MC can be seen in the book, table 7.1.

Theorem 1.12 (Convergence of TD)

Given a policy π , in 1.23, $v_t(s)$ converges almost surely to $V^\pi(s)$ as $t \rightarrow \infty$ for all $s \in \mathcal{S}$ if $\sum_t \alpha_t(s) = \infty$ and $\sum_t \alpha_t^2(s) < \infty$ for all $s \in \mathcal{S}$. When α_t is constant, it can still be shown that the algorithm converges in the sense of expectation, though it does not follow second condition.



Proof Deducting $V^\pi(s)$ from both side in 1.23, then show that it follows 1.10.

Then we introduce *Sarsa*, another TD algorithm that can directly estimate action values.

$$q_{t+1}(s_t, a_t) = q_t(s_t, a_t) - \alpha_t(s_t, a_t)[q_t(s_t, a_t) - (r_{t+1} + \gamma q_t(s_{t+1}, a_{t+1}))] \quad (1.24)$$

and $q_{t+1}(s, a) = q_t(s, a)$ for all $(s, a) \neq (s_t, a_t)$, where $t = 0, 1, 2, \dots$. It is given for solving the Bellman equation of a given policy:

$$Q^\pi(s, a) = \mathbb{E}[R + \gamma Q^\pi(S', A') \mid s, a], \quad \forall (s, a) \quad (1.25)$$

Let's proof this is a Bellman equation:

Proof

$$\begin{aligned} Q^\pi(s, a) &= \sum_r r p(r \mid s, a) + \gamma \sum_{s'} \sum_{a'} Q^\pi(s', a') p(s' \mid s, a) \pi(a' \mid s') \\ &= \sum_r r p(r \mid s, a) + \gamma \sum_{s'} p(s' \mid s, a) \sum_{a'} Q^\pi(s', a') \pi(a' \mid s'). \end{aligned}$$

This equation establishes the relationships among the action values. Since

$$\begin{aligned} p(s', a' \mid s, a) &= p(s' \mid s, a) p(a' \mid s', s, a) \\ &= p(s' \mid s, a) p(a' \mid s') \quad \text{due to conditional independence} \\ &= p(s' \mid s, a) \pi(a' \mid s') \end{aligned}$$

Then

$$Q^\pi(s, a) = \sum_r rp(r|s, a) + \gamma \sum_{s'} \sum_{a'} Q^\pi(s', a') p(s', a'|s, a).$$

Algorithm 3: Optimal policy learning by Sarsa

Input: Initial policy $\pi_0(a|s)$, $Q_0(s, a)$, $\forall(s, a)$, $\epsilon \in (0, 1]$

```

1 foreach episode do
2   Episode generation: Generate  $a_0$  at  $s_0$  following  $\pi_0(s_0)$ 
3   while  $s_t(t = 0, 1, 2, \dots) \neq s_{target}$  do
4     Generate  $r_{t+1}, s_{t+1}$  and  $a_{t+1}$ , generate  $a_{t+1}$  following  $\pi_{t+1}(s_t)$ , generate  $r_{t+1}, s_{t+1}$  by interacting with
       environment
5      $q_{t+1}(s_t, a_t) \leftarrow q_t(s_t, a_t) - \alpha_t(s_t, a_t)[q_t(s_t, a_t) - (r_{t+1} + \gamma q_t(s_{t+1}, a_{t+1}))]$ 
6      $\pi_{t+1}(a|s_t) \leftarrow \begin{cases} 1 - \frac{\epsilon}{|\mathcal{A}(s)|} (|\mathcal{A}(s)| - 1) & a = a^* = \arg \max_a q_{t+1}(s_t, a) \\ \frac{\epsilon}{|\mathcal{A}(s)|} & a \neq a^*. \end{cases}$ 
7   end
8 end
```

A variant of Sarsa is Expected Sarsa:

$$q_{t+1}(s_t, a_t) = q_t(s_t, a_t) - \alpha_t(s_t, a_t)[q_t(s_t, a_t) - (r_{t+1} + \gamma \mathbb{E}[q_t(s_{t+1}, A)])]$$

where $\mathbb{E}[q_t(s_{t+1}, A)] = \sum_a \pi_t(a|s_{t+1}) q_t(s_{t+1}, a) = v_t(s_{t+1})$. Although calculating the expected value may increase the computational complexity slightly, it is beneficial in the sense that it reduces the estimation variances. This variant can be viewed as solving

$$Q^\pi(s, a) = \mathbb{E}[R_{t+1} + \gamma \mathbb{E}[Q^\pi(S_{t+1}, A_{t+1}) | S_{t+1} | S_t = s, A_t = a]], \quad \forall(s, a)$$

which is again a Bellman equation as $\mathbb{E}[Q^\pi(S_{t+1}, A_{t+1}) | S_{t+1}] = \sum_{A'} Q^\pi(S_{t+1}, A') \pi(A' | S_{t+1}) = V^\pi(S_{t+1})$.

Now we introduce **n-step Sarsa**, We will see that Sarsa and MC learning are two extreme cases of n-step Sarsa.

$$\begin{aligned}
\text{Sarsa} &\leftarrow \begin{aligned} G_t^{(1)} &= R_{t+1} + \gamma Q^\pi(S_{t+1}, A_{t+1}) \\ G_t^{(2)} &= R_{t+1} + \gamma R_{t+2} + \gamma^2 Q^\pi(S_{t+2}, A_{t+2}) \\ &\vdots \end{aligned} \\
\text{n-step Sarsa} &\leftarrow \begin{aligned} G_t^{(n)} &= R_{t+1} + \gamma R_{t+2} + \dots + \gamma^n Q^\pi(S_{t+n}, A_{t+n}) \\ &\vdots \end{aligned} \\
\text{policy iteration} &\leftarrow G_t^{(\infty)} = R_{t+1} + \gamma R_{t+2} + \dots
\end{aligned}$$

It should be noted that $G_t = G_t^{(1)} = G_t^{(2)} = G_t^{(n)} = G_t^{(\infty)}$, where the superscripts merely indicate the different decomposition structures of G_t .

When $n = 1$, we get 1.24, when $n = \infty$, $q_{t+1}(s_t, a_t) = g_t = r_{t+1} + \gamma r_{t+2} + \gamma^2 r_{t+3} + \dots$. And for the general n :

$$q_{t+1}(s_t, a_t) = q_t(s_t, a_t) - \alpha_t(s_t, a_t)[q_t(s_t, a_t) - (r_{t+1} + \gamma r_{t+2} + \dots + \gamma^n q_t(s_{t+n}, a_{t+n}))]$$

Now, finally we introduce **Q-learning**, which can directly estimate optimal action values and find optimal policies (Sarsa need the policy improvement step):

$$q_{t+1}(s_t, a_t) = q_t(s_t, a_t) - \alpha_t(s_t, a_t)[q_t(s_t, a_t) - (r_{t+1} + \gamma \max_{a \in A(s_{t+1})} q_t(s_{t+1}, a))] \quad (1.26)$$


Given (s_t, a_t) , Sarsa need $(r_{t+1}, s_{t+1}, a_{t+1})$ in every iteration, whereas Q-learning merely requires (r_{t+1}, s_{t+1}) . It is for solving the following equation:

$$q(s, a) = \mathbb{E}[R_{t+1} + \gamma \max_a q(S_{t+1}, a) | S_t = s, A_t = a] \quad (1.27)$$

we can show this is also a Bellman equation:

$$\begin{aligned}
 q(s, a) &= \sum_r p(r | s, a) r + \gamma \sum_{s'} p(s' | s, a) \max_{a \in \mathcal{A}(s')} q(s', a) \\
 &\xrightarrow{\text{taking maximum}} \\
 \max_{a \in \mathcal{A}(s)} q(s, a) &= \max_{a \in \mathcal{A}(s')} \left[\sum_r p(r | s, a) r + \gamma \sum_{s'} p(s' | s, a) \max_{a \in \mathcal{A}(s')} q(s', a) \right] \\
 &\xrightarrow{\text{let } v(s) = \max_{a \in \mathcal{A}(s)} q(s, a)} \\
 v(s) &= \max_{a \in \mathcal{A}(s')} \left[\sum_r p(r | s, a) r + \gamma \sum_{s'} p(s' | s, a) v(s') \right] \\
 &= \max_{\pi} \sum_{a \in \mathcal{A}(s)} \pi(a | s) \left[\sum_r p(r | s, a) r + \gamma \sum_{s'} p(s' | s, a) v(s') \right]
 \end{aligned}$$

Definition 1.3

- **behavior policy** is the one used to generate experience samples.
- **target policy** is the one that is constantly updated to converge to an optimal policy.
- **on-policy**: When the behavior policy is the same as the target policy, such a learning process is called on-policy.
- **off-policy**: When the behavior policy is **not** the same as the target policy, like Q-learning. The advantage of off-policy learning is that it can learn optimal policies based on the experience samples generated by other policies. So we can choose a exploratory behavior policy and the learning efficiency will be increased. 

The fundamental reason why Q-learning is off-policy is that Q-learning aims to solve the Bellman optimality equation rather than the Bellman equation of a given policy.

Algorithm 4: Optimal policy learning via Q-learning (on-policy version)

Input: Initial policy $\pi_0(a|s)$ from $q_0, q_0(s, a), \forall(s, a), \epsilon \in (0, 1]$

Output: Learn an optimal path that can lead the agent to the target state from an initial state s_0

```

1 foreach episode do
2   while  $s_t (t = 0, 1, 2, \dots) \neq s_{target}$  do
3     Collect  $a_t, r_{t+1}, s_{t+1}$ , generate  $a_t$  following  $\pi_t(s_t)$ , generate  $r_{t+1}, s_{t+1}$  by interacting with enviroment
4      $q_{t+1}(s_t, a_t) = q_t(s_t, a_t) - \alpha_t(s_t, a_t) [q_t(s_t, a_t) - (r_{t+1} + \gamma \max_{a \in \mathcal{A}(s_{t+1})} q_t(s_{t+1}, a))]$ 
5      $\pi_{t+1}(a|s_t) \leftarrow \begin{cases} 1 - \frac{\epsilon}{|\mathcal{A}(s)|} (|\mathcal{A}(s)| - 1) & a = a^* = \arg \max_a q_{t+1}(s_t, a) \\ \frac{\epsilon}{|\mathcal{A}(s)|} & a \neq a^*. \end{cases}$ 
6   end
7 end

```

Algorithm 5: Optimal policy learning via Q-learning (off-policy version)

Input: Behavior policy $\pi_b(a|s)$, initial guess $q_0(s, a), \forall(s, a)$ and all t

Output: Learn an optimal target policy π_T for all states from the experience samples generated by π_b

```

1 foreach episode  $\{s_0, a_0, r_1, s_1, a_1, r_2, \dots\}$  generated by  $\pi_b$  do
2   foreach step  $t = 0, 1, 2, \dots$  of the episode do
3      $q_{t+1}(s_t, a_t) = q_t(s_t, a_t) - \alpha_t(s_t, a_t) [q_t(s_t, a_t) - (r_{t+1} + \gamma \max_{a \in \mathcal{A}(s_{t+1})} q_t(s_{t+1}, a))]$ 
4      $\pi_{T,t+1}(a|s_t) \leftarrow \begin{cases} 1 & a = a^* = \arg \max_a q_{t+1}(s_t, a) \\ 0 & a \neq a^*. \end{cases}$ 
5   end
6 end

```

We can use a unified formulation for all the TD algorithm:

$$q_{t+1}(s_t, a_t) = q_t(s_t, a_t) - \alpha_t(s_t, a_t)[q_t(s_t, a_t) - \bar{q}_t]$$

Algorithm	Expression of the TD target \bar{q}_t in (7.20)
Sarsa	$\bar{q}_t = r_{t+1} + \gamma q_t(s_{t+1}, a_{t+1})$
n -step Sarsa	$\bar{q}_t = r_{t+1} + \gamma r_{t+2} + \dots + \gamma^n q_t(s_{t+n}, a_{t+n})$
Q-learning	$\bar{q}_t = r_{t+1} + \gamma \max_a q_t(s_{t+1}, a)$
Monte Carlo	$\bar{q}_t = r_{t+1} + \gamma r_{t+2} + \gamma^2 r_{t+3} + \dots$

Algorithm	Equation to be solved
Sarsa	BE: $q_\pi(s, a) = \mathbb{E}[R_{t+1} + \gamma q_\pi(S_{t+1}, A_{t+1}) S_t = s, A_t = a]$
n -step Sarsa	BE: $q_\pi(s, a) = \mathbb{E}[R_{t+1} + \gamma R_{t+2} + \dots + \gamma^n q_\pi(S_{t+n}, A_{t+n}) S_t = s, A_t = a]$
Q-learning	BOE: $q(s, a) = \mathbb{E}[R_{t+1} + \gamma \max_a q(S_{t+1}, a) S_t = s, A_t = a]$
Monte Carlo	BE: $q_\pi(s, a) = \mathbb{E}[R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots S_t = s, A_t = a]$

Figure 1.2: TD algorithms

1.11 Value Function Methods

Instead of using table to represent the value function, which takes a lot of memory and does not generalize well at unknown state, we now introduce another way - using function estimation. In particular, the objective function is:

$$J(w) = \mathbb{E}[(V^\pi(S) - \hat{v}(S, w))^2] \quad (1.28)$$

We can use a uniform distribution as the distribution of S , but it does not consider the Markov process property. So we use **stationary distribution**, the long-term behavior, after sufficiently long period, every state will have its stationary probability. Let $\{d_\pi(s)\}_{s \in \mathcal{S}}$ denote the stationary distribution of the Markov process under policy π . The objective is now:

$$J(w) = \sum_{s \in \mathcal{S}} d_\pi(s) (V^\pi(S) - \hat{v}(S, w))^2 \quad (1.29)$$

Definition 1.4 (Stationary Distribution)

The probability of the agent transitioning from s_i to s_j using exactly k steps is denoted as

$$p_{ij}^k = \Pr(S_{t_k} = j \mid S_{t_0} = i)$$

by the definition of P_π , we have

$$[P_\pi]_{ij} = p_{ij}^{(1)}$$

which means $[P_\pi]_{ij}$ is the probability of transitioning from s_i to s_j using a single step. Second, consider P_π^2 :

$$[P_\pi^2]_{ij} = [P_\pi P_\pi]_{ij} = \sum_{q=1}^n [P_\pi]_{iq} [P_\pi]_{qj} = p_{ij}^{(2)}$$

In the same manner:

$$[P_\pi^n]_{ij} = p_{ij}^{(n)}$$

Let $d_0 \in \mathcal{R}^n$ be a vector representing the probability distribution of the states at the initial time step. Then

$$d_k(s_i) = \sum_{j=1}^n d_0(s_j) [P_\pi^k]_{ji}, i = 1, 2, \dots \quad (1.30)$$

The matrix form is:

$$d_k^T = d_0^T P_\pi^k$$

When we consider the long-term behavior of the Markov process, it holds under certain conditions that

$$\lim_{k \rightarrow \infty} P_\pi^k = \mathbf{1}_n d_\pi^T$$

$\mathbf{1}_n d_\pi^T$ is a constant matrix with all its rows equal to d_π^T , so

$$\lim_{k \rightarrow \infty} d_k^T = d_0^T \lim_{k \rightarrow \infty} P_\pi^k = d_0^T \mathbf{1}_n d_\pi^T = d_\pi^T \quad (1.31)$$

this means that the state distribution converges to a constant value d_π , which is called the limiting distribution, it is independent of initial distribution d_0 .

The value of d_π is by taking limit of $d_k^T = d_{k-1}^T P_\pi$, where we will get

$$d_\pi^T = d_\pi^T P_\pi \quad (1.32)$$

As a result, it is the left eigenvector of P_π associated with eigenvalue 1.

The solution of 1.32 is usually called a stationary distribution, whereas the distribution in 1.31 is usually called the limiting distribution. Note that 1.31 implies 1.32, but the converse may not be true.

- State s_j is said to be **accessible** from s_i if there exists a finite integer k so that $[P_\pi^k]_{ij} > 0$.
- If two states s_i and s_j are mutually accessible, then the two states are said to **communicate**.
- A Markov process is called **irreducible** if all of its states communicate with each other. k maybe different for different i, j .
- A Markov process is called **regular** if there exists $k \geq 1$, such that $[P_\pi^k]_{ij} > 0, \forall i, j$. Regular Markov process is also irreducible, but the inverse is not. However, if a irreducible one exists i such that $[P_\pi]_{ij} > 0$, then it is also regular. Moreover, if $P_\pi^k > 0$, then for any $k \geq k', P_\pi^{k'} > 0$.



We can use gradient descent to minimize the $J(w)$:

$$\begin{aligned} \nabla_w J(w_k) &= \nabla_w \mathbb{E}[(v_\pi(S) - \hat{v}(S, w_k))^2] \\ &= \mathbb{E}[\nabla_w (v_\pi(S) - \hat{v}(S, w_k))^2] \\ &= 2\mathbb{E}[(v_\pi(S) - \hat{v}(S, w_k))(-\nabla_w \hat{v}(S, w_k))] \\ &= -2\mathbb{E}[(v_\pi(S) - \hat{v}(S, w_k))\nabla_w \hat{v}(S, w_k)] \end{aligned}$$

therefore the algorithm is:

$$w_{k+1} = w_k + 2\alpha_k \mathbb{E}[(v_\pi(S) - \hat{v}(S, w_k))\nabla_w \hat{v}(S, w_k)] \quad (1.33)$$

with the stochastic replacement, and coefficient 2 merge into α_k without loss of generality, we get:

$$w_{t+1} = w_t + \alpha_t (v_\pi(s_t) - \hat{v}(s_t, w_t)) \nabla_w \hat{v}(s_t, w_t) \quad (1.34)$$

Now we need to replace v_π , either with Monte Carlo:

$$w_{t+1} = w_t + \alpha_t (g_t - \hat{v}(s_t, w_t)) \nabla_w \hat{v}(s_t, w_t)$$

or with TD

$$w_{t+1} = w_t + \alpha_t [r_{t+1} + \gamma \hat{v}(s_{t+1}, w_t) - \hat{v}(s_t, w_t)] \nabla_w \hat{v}(s_t, w_t) \quad (1.35)$$

Tabular TD learning is a special case of TD-Linear.

However, 1.35 does not really solve 1.28, so we need to explain it rigorously. We first consider deterministic algorithm:

$$w_{t+1} = w_t + \alpha_t \mathbb{E} \left[(r_{t+1} + \gamma \phi^T(s_{t+1}) w_t - \phi^T(s_t) w_t) \phi(s_t) \right] \quad (1.36)$$

then we will show 1.35 is the SGD implementation of it. Let

$$\Phi = \begin{bmatrix} \vdots \\ \phi^T(s) \\ \vdots \end{bmatrix} \in \mathbb{R}^{n \times m}, \quad D = \begin{bmatrix} \ddots & & \\ & d_\pi(s) & \\ & & \ddots \end{bmatrix} \in \mathbb{R}^{n \times n}$$

Lemma 1.2

The expectation in 1.35 can be rewritten as

$$\mathbb{E} \left[(r_{t+1} + \gamma \phi^T(s_{t+1})w_t - \phi^T(s_t)w_t) \phi(s_t) \right] = b - Aw_t$$

where

$$A \doteq \Phi^T D (I - \gamma P_\pi) \Phi \in \mathbb{R}^{m \times m}$$

$$b \doteq \Phi^T D r_\pi \in \mathbb{R}^m.$$

and s_t is assumed to obey the stationary distribution d_π .



Then 1.35 can be rewritten as $w_{t+1} = w_t + \alpha_t(b - Aw_t)$, and it converges to $w^* = A^{-1}b$, since A is invertible and positive definite. Moreover, in the tabular case, as $n = |S|$ and $\phi(s) = [0, \dots, 1, \dots, 0]^T$, we will get $w^* = A^{-1}b = v_\pi$. These statements' proof please refer to the book section 8.2.

Then we show TD learning algorithm 1.35 minimizes the projected Bellman error. 1.28 can be written in vector form as $J(w) = \|\hat{v}(w) - v_\pi\|_D^2$, where $\|x\|_D^2 = x^T D x = \|D^{1/2}x\|_2^2$, D is given above. But this one has the true state value, so now we introduce Bellman error

$$J_{BE}(w) = \|\hat{v}(w) - (r_\pi + \gamma P_\pi \hat{v}(w))\|_D^2 = \|\hat{v}(w) - T_\pi(\hat{v}(w))\|_D^2 \quad (1.37)$$

where $T_\pi(x) = r_\pi + \gamma P_\pi x$ is the Bellman operator.

We can see this error may not be minimized to zero due to the limited approximation ability of the approximator. By contrast, an objective function that can be minimized to zero is the projected Bellman error

$$J_{PBE}(w) = \|\hat{v}(w) - MT_\pi(\hat{v}(w))\|_D^2 \quad (1.38)$$

where $M \in \mathbb{R}^{n \times n}$ is the orthogonal projection matrix that geometrically projects any vector onto the space of all approximations. We can show that $M = \Phi(\Phi^T D \Phi)^{-1} \Phi^T D \in \mathbb{R}^{n \times n}$ is the projection matrix that geometrically projects any vector onto the range space of Φ . Since $\hat{v}(w)$ is in the range space of Φ , we can always find a value of w that can minimize $J_{PBE}(w)$ to zero, it can be proven it is $w^* = A^{-1}b$.

We can show that in linear case the bound of the error

$$\|\hat{v}(w^*) - v_\pi\|_D = \|\Phi w^* - v_\pi\|_D \leq \frac{1}{1 - \gamma} \min_w \|\hat{v}(w) - v_\pi\|_D = \frac{1}{1 - \gamma} \min_w \sqrt{J_E(w)}$$

Next we introduce *least-squares* TD (LSTD), we can write A and b as (in proof of 1.2)

$$A = \mathbb{E} \left[\phi(s_t) (\phi(s_t) - \gamma \phi(s_{t+1}))^T \right],$$

$$b = \mathbb{E} [r_{t+1} \phi(s_t)].$$

The idea of LSTD is simple: if we can use random samples to directly obtain the estimates of A and b , then the optimal parameter can be directly estimated. In particular, given a trajectory (s_0, r_1, s_1, \dots) under π , we have

$$\hat{A}_t = \sum_{k=0}^{t-1} \phi(s_k) (\phi(s_k) - \gamma \phi(s_{k+1}))^T, \quad (1.39)$$

$$\hat{b}_t = \sum_{k=0}^{t-1} r_{k+1} \phi(s_k). \quad (1.40)$$

we omitted $1/t$, because it does not change the result that $w_t = \hat{A}_t^{-1} \hat{b}_t$.

The advantage of LSTD is that it uses experience samples more efficiently and converges faster than the TD method.

The disadvantages of LSTD are as follows. First, it can only estimate state values. Moreover, while the TD algorithm allows nonlinear approximators, LSTD does not. Second, the computational cost of LSTD is higher than that of TD since

LSTD as we need to calculate inverse (but we can introduce a iterative algorithm).

Now, all we talked in this section evaluate value function, but what about action values? We can modify 1.35

$$w_{t+1} = w_t + \alpha_t \left[r_{t+1} + \gamma \hat{q}(s_{t+1}, a_{t+1}, w_t) - \hat{q}(s_t, a_t, w_t) \right] \nabla_w \hat{q}(s_t, a_t, w_t) \quad (1.41)$$

Algorithm 6: Sarsa with function approximation

Input: Initial policy π_0 , initial parameter w_0

Output: Learn an optimal policy that can lead the agent to the target state from an initial state s_0

```

1 foreach episode do
2   Generate  $a_0$  at  $s_0$  following  $\pi_0(s_0)$ 
3   while  $s_t (t = 0, 1, 2, \dots) \neq s_{target}$  do
4     Collect  $(r_{t+1}, s_{t+1}, a_{t+1})$  given  $(s_t, a_t)$ , last one by following  $\pi_t(s_{t+1})$ , other by interacting with
       enviroment
5      $w_{t+1} = w_t + \alpha_t \left[ r_{t+1} + \gamma \hat{q}(s_{t+1}, a_{t+1}, w_t) - \hat{q}(s_t, a_t, w_t) \right] \nabla_w \hat{q}(s_t, a_t, w_t)$ 
6      $\pi_{t+1}(a|s_t) \leftarrow \begin{cases} 1 - \frac{\epsilon}{|\mathcal{A}(s)|} (|\mathcal{A}(s)| - 1) & a = a^* = \arg \max_a q_{t+1}(s_t, a) \\ \frac{\epsilon}{|\mathcal{A}(s)|} & a \neq a^*. \end{cases}$ 
7      $s_t \leftarrow s_{t+1}, a_t \leftarrow a_{t+1}$ 
8   end
9 end

```

We can also modify Q-learning

$$w_{t+1} = w_t + \alpha_t \left[r_{t+1} + \gamma \max_{a \in A(s_{t+1})} \hat{q}(s_{t+1}, a, w_t) - \hat{q}(s_t, a_t, w_t) \right] \nabla_w \hat{q}(s_t, a_t, w_t) \quad (1.42)$$

again we can implement a on-policy or off-policy version.

Lastly, in this section we introduce **Deep Q-learning** (DQN), the objective function is

$$J = \mathbb{E} \left[\left(R + \gamma \max_{a \in A(S')} \hat{q}(S', a, w) - \hat{q}(S, A, w) \right)^2 \right] \quad (1.43)$$

This objective function can be viewed as the squared Bellman optimality error. That is because

$$q(s, a) = \mathbb{E} \left[R_{t+1} + \gamma \max_{a \in A(S_{t+1})} q(S_{t+1}, a) \mid S_t = s, A_t = a \right]$$

is the Bellman optimality equation (see last section).

For the sake of simplicity, it is assumed that the value of w in $\gamma \max_{a \in A(S_{t+1})} q(S_{t+1}, a)$ is fixed (for a short period of time) so the calculation of the gradient becomes much easier. Let

$$J = \mathbb{E} \left[\left(R + \gamma \max_{a \in A(S')} \hat{q}(S', a, w_T) - \hat{q}(S, A, w) \right)^2 \right]$$

where w_T is the target network's parameter which is fixed. Then the gradient is:

$$\nabla_w J = -\mathbb{E} \left[\left(R + \gamma \max_{a \in A(S')} \hat{q}(S', a, w_T) - \hat{q}(S, A, w) \right) \nabla_w \hat{q}(S, A, w) \right] \quad (1.44)$$

We need to pay attention that:

- The first technique is to use two networks, a main network and a target network, as mentioned when we calculate the gradient, let w and w_T denote the parameters of the main and target networks, respectively. They are initially set to the same value. The main network is updated in every iteration. By contrast, the target network is set to be the same as the main network every certain number of iterations to satisfy the assumption that w_T is fixed when calculating the gradient.
- One notable difference between deep and nondeep reinforcement learning algorithms is that a mini-batch of samples to used to train a network in DQN instead of using a single sample to update the main network.
- Let (s, a, r, s') be an experience sample and $\mathcal{B} \doteq \{(s, a, r, s')\}$ the **replay buffer**. Every time we update the main

network, we can draw a mini-batch of experience samples from the replay buffer. The draw of samples, or called experience replay, **should follow a uniform distribution**. The state-action samples may not be uniformly distributed in practice since they are generated as a sequence by the behavior policy. It is necessary to break the correlation between the samples in the sequence to satisfy the assumption of uniform distribution. To do this, we can use the experience replay technique by uniformly drawing samples from the replay buffer. A benefit of experience replay is that each experience sample may be used multiple times, which can increase the data efficiency.

Algorithm 7: Deep Q-learning (off-policy version)

Input: A main network and a target network with the same initial parameter.

Output: Learn an optimal target network to approximate the optimal action values from the experience samples generated by a given behavior policy π_b

```

1 begin
2   Store the experience samples generated by  $\pi_b$  in a replay buffer  $\mathcal{B} \doteq \{(s, a, r, s')\}$ 
3   foreach iteration do
4     Uniformly draw a mini-batch of samples from  $\mathcal{B}$ 
5     foreach sample  $(s_i, a_i, r_i, s'_i)$  do
6        $y_i \doteq r_i + \gamma \max_{a \in A(s'_i)} \hat{q}(s'_i, a, w_T)$ 
7     end
8     Update the main network to minimize  $(y_T - \hat{q}(s, a, w))^2$  using the mini-batch of samples, for example,
        $L(w) = \frac{1}{n} \sum_{i=1}^N (y_i - \hat{q}(s_i, a_i, w))^2$ .
9   end
10  Set  $w_T = w$  every  $C$  iteration
11 end

```

1.12 Policy Gradient Methods

First, how to define optimal policies? When represented as a table, a policy is defined as optimal if it can maximize **every state value**. When represented by a function, a policy is defined as optimal if it can maximize certain **scalar metrics**.

We now introduce some metrics

Definition 1.5 (Average State Value)

Or average value, is defined as $\bar{v}_\pi = \sum_{s \in \mathcal{S}} d(s) v_\pi(s)$, where $d(s)$ is the weight of the state, It satisfies $d(s) \geq 0$ for any $s \in \mathcal{S}$ and $\sum_{s \in \mathcal{S}} d(s) = 1$. So we can interpret it as the probability distribution, so $\bar{v}_\pi = \mathbb{E}_{S \sim d}[v_\pi(S)]$. Then how can we select d , we can choose it independent (we use d_0 and \bar{v}_π^0 as the notation) or dependent of policy (the stationary distribution). Our ultimate goal is to find an optimal policy to maximize \bar{v}_π . It is the same as

$$J(\theta) = \lim_{n \rightarrow \infty} \mathbb{E} \left[\sum_{t=0}^n \gamma^t R_{t+1} \right] = \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t R_{t+1} \right]$$

because

$$\begin{aligned}
 \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t R_{t+1} \right] &= \sum_{s \in \mathcal{S}} d(s) \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t R_{t+1} \mid S_0 = s \right] \\
 &= \sum_{s \in \mathcal{S}} d(s) v_\pi(s) \\
 &= \bar{v}_\pi(s)
 \end{aligned}$$

with the matrix form we have

$$\bar{v}_\pi(s) = d^T v_\pi(s)$$



Definition 1.6 (Average Reward)

Or called average one-step reward, it is defined as

$$\begin{aligned}\bar{r}_\pi &\doteq \sum_{s \in \mathcal{S}} d_\pi(s) r_\pi(s) \\ &= \mathbb{E}_{S \sim d_\pi} [r_\pi(S)]\end{aligned}$$

where $r_\pi(s) \doteq \sum_{a \in \mathcal{A}} \pi(a|s, \theta) r(s, a) = \mathbb{E}_{A \in \pi(s, \theta)} [r(s, A)|s]$ is the expectation of the immediate rewards. Here $r(s, a) \doteq \mathbb{E}[R|s, a] = \sum_r r p(r|s, a)$.

\bar{r}_π is equivalent to

$$J(\theta) = \lim_{n \rightarrow \infty} \frac{1}{n} \mathbb{E} \left[\sum_{t=0}^{n-1} R_{t+1} \right]$$

and

$$\bar{r}_\pi = \sum_{s \in \mathcal{S}} d_\pi(s) r_\pi(s) = d_\pi^T r_\pi$$



Table 1.1: Summary of the different but equivalent expressions of \bar{v}_π and \bar{r}_π

Metric	Expression 1	Expression 2	Expression 3
\bar{v}_π	$\sum_{s \in \mathcal{S}} d(s) v_\pi(s)$	$\mathbb{E}_{S \sim d} [v_\pi(S)]$	$\lim_{n \rightarrow \infty} \mathbb{E} \left[\sum_{t=0}^n \gamma^t R_{t+1} \right]$
\bar{r}_π	$\sum_{s \in \mathcal{S}} d_\pi(s) r_\pi(s)$	$\mathbb{E}_{S \sim d_\pi} [r_\pi(S)]$	$\lim_{n \rightarrow \infty} \frac{1}{n} \mathbb{E} \left[\sum_{t=0}^{n-1} R_{t+1} \right]$

Theorem 1.13 (Policy gradient theorem)

The gradient of $J(\theta)$ is

$$\nabla_\theta J(\theta) = \sum_{s \in \mathcal{S}} \eta(s) \sum_{a \in \mathcal{A}} \nabla_\theta \pi(a|s, \theta) q_\pi(s, a) \quad (1.45)$$

where η is a state distribution and $\nabla_\theta \pi$ is the gradient of π with respect to θ . It has a compact form

$$\nabla_\theta J(\theta) = \mathbb{E}_{S \sim \eta, A \sim \pi(S, \theta)} [\nabla_\theta \ln \pi(A|S, \theta) q_\pi(S, A)] \quad (1.46)$$

we can use softmax to make sure $\pi(a|s, \theta) \in (0, 1), \forall s \in \mathcal{S}$

$$\pi(a|s, \theta) = \frac{e^{h(s, a, \theta)}}{\sum_{a' \in \mathcal{A}} e^{h(s, a', \theta)}}, \quad a \in \mathcal{A} \quad (1.47)$$

and in this sense the policy is exploratory.

**Lemma 1.3**

$$\bar{r}_\pi = (1 - \gamma) \bar{v}_\pi \quad (1.48)$$



Proof Since $v_\pi = r_\pi + \gamma P_\pi v_\pi$, multiply d_π^T on both side, we get

$$\bar{v}_\pi = \bar{r}_\pi + \gamma d_\pi^T P_\pi v_\pi = \bar{r}_\pi + \gamma d_\pi^T v_\pi = \bar{r}_\pi + \gamma \bar{v}_\pi$$

Lemma 1.4

In the discounted case, it holds for any $s \in \mathcal{S}$ that

$$\nabla_\theta v_\pi(s) = \sum_{s' \in \mathcal{S}} \Pr_\pi(s'|s) \sum_{a \in \mathcal{A}} \nabla_\theta \pi(a|s', \theta) q_\pi(s', a) \quad (1.49)$$

where

$$\Pr_\pi(s'|s) \doteq \sum_{k=0}^{\infty} \gamma^k [P_\pi^k]_{ss'} = [(I_n - \gamma P_\pi)^{-1}]_{ss'}$$

is the discounted total probability of transitioning from s to s' under policy π . Here $[\cdot]_{ss'}$ denotes the entry in the s th row and s' th column, k means exactly k steps.

Proof

$$\begin{aligned}\nabla_{\theta} v_{\pi}(s) &= \nabla_{\theta} \left[\sum_{a \in A} \pi(a|s, \theta) q_{\pi}(s, a) \right] \\ &= \sum_{a \in A} \left[\nabla_{\theta} \pi(a|s, \theta) q_{\pi}(s, a) + \pi(a|s, \theta) \nabla_{\theta} q_{\pi}(s, a) \right]\end{aligned}\quad (1.50)$$

where

$$q_{\pi}(s, a) = r(s, a) + \gamma \sum_{s' \in \mathcal{S}} p(s'|s, a) v_{\pi}(s')$$

since $r(s, a) = \sum_r r p(r|s, a)$ is independent of θ

$$\nabla_{\theta} q_{\pi}(s, a) = 0 + \gamma \sum_{s' \in \mathcal{S}} p(s'|s, a) \nabla_{\theta} v_{\pi}(s')$$

substitute this into the first equation we have

$$\begin{aligned}\nabla_{\theta} v_{\pi}(s) &= \sum_{a \in A} \left[\nabla_{\theta} \pi(a|s, \theta) q_{\pi}(s, a) + \pi(a|s, \theta) \gamma \sum_{s' \in \mathcal{S}} p(s'|s, a) \nabla_{\theta} v_{\pi}(s') \right] \\ &= \sum_{a \in A} \nabla_{\theta} \pi(a|s, \theta) q_{\pi}(s, a) + \gamma \sum_{a \in A} \pi(a|s, \theta) \sum_{s' \in \mathcal{S}} p(s'|s, a) \nabla_{\theta} v_{\pi}(s')\end{aligned}\quad (1.51)$$

One way to calculate it is to use the **unrolling technique**, we use another way based on the matrix-vector form, let

$$u(s) \doteq \sum_{a \in A} \nabla_{\theta} \pi(a|s, \theta) q_{\pi}(s, a)$$

Since

$$\sum_{a \in A} \pi(a|s, \theta) \sum_{s' \in \mathcal{S}} p(s'|s, a) \nabla_{\theta} v_{\pi}(s') = \sum_{s' \in \mathcal{S}} p(s'|s) \nabla_{\theta} v_{\pi}(s') = \sum_{s' \in \mathcal{S}} [P_{\pi}]_{ss'} \nabla_{\theta} v_{\pi}(s')$$

the matrix form is

$$\underbrace{\begin{bmatrix} \vdots \\ \nabla_{\theta} v_{\pi}(s) \\ \vdots \end{bmatrix}}_{\nabla_{\theta} v_{\pi} \in \mathbb{R}^{mn}} = \underbrace{\begin{bmatrix} \vdots \\ u(s) \\ \vdots \end{bmatrix}}_{u \in \mathbb{R}^{mn}} + \gamma (P_{\pi} \otimes I_m) \underbrace{\begin{bmatrix} \vdots \\ \nabla_{\theta} v_{\pi}(s') \\ \vdots \end{bmatrix}}_{\nabla_{\theta} v_{\pi} \in \mathbb{R}^{mn}}$$

which can be written concisely as

$$\nabla_{\theta} v_{\pi} = u + \gamma (P_{\pi} \otimes I_m) \nabla_{\theta} v_{\pi}$$

here $n = |\mathcal{S}|$ and m is the dimension of the parameter vector θ . Then

$$\begin{aligned}\nabla_{\theta} v_{\pi} &= (I_{nm} - \gamma P_{\pi} \otimes I_m)^{-1} u \\ &= (I_n \otimes I_m - \gamma P_{\pi} \otimes I_m)^{-1} u \\ &= [(I_n - \gamma P_{\pi})^{-1} \otimes I_m] u\end{aligned}\quad (1.52)$$

from this we get

$$\begin{aligned}\nabla_{\theta} v_{\pi}(s) &= \sum_{s' \in \mathcal{S}} [(I_n - \gamma P_{\pi})^{-1}]_{ss'} u(s') \\ &= \sum_{s' \in \mathcal{S}} [(I_n - \gamma P_{\pi})^{-1}]_{ss'} \sum_{a \in A} \nabla_{\theta} \pi(a|s', \theta) q_{\pi}(s', a)\end{aligned}\quad (1.53)$$

since $(I_n - \gamma P_{\pi})^{-1} = I + \gamma P_{\pi} + \gamma^2 P_{\pi}^2 + \dots$

$$[(I_n - \gamma P_{\pi})^{-1}]_{ss'} = [I]_{ss'} + \gamma [P_{\pi}]_{ss'} + \gamma^2 [P_{\pi}^2]_{ss'} + \dots = \sum_{k=0}^{\infty} \gamma^k [P_{\pi}^k]_{ss'}$$

denote $[(I_n - \gamma P_{\pi})^{-1}]_{ss'} \doteq \Pr_{\pi}(s'|s)$, we get 1.49.

Theorem 1.14 (Gradient of \bar{v}_{π}^0 in the discounted case)

In the discounted case where $\gamma \in (0, 1)$, the gradient of \bar{v}_{π}^0 is

$$\nabla_{\theta} \bar{v}_{\pi}^0 = \mathbb{E}[\nabla_{\theta} \ln \pi(A|S, \theta) q_{\pi}(S, A)]$$

where $S \sim \rho_\pi$ and $A \sim \pi(S, \theta)$, the state distribution is

$$\rho_\pi(s) = \sum_{s' \in \mathcal{S}} d_0(s') \Pr_\pi(s|s'), \quad s \in \mathcal{S} \quad (1.54)$$

Proof Since $d_0(s)$ is independent of π , we have

$$\nabla_\theta \bar{v}_\pi^0 = \nabla_\theta \sum_{s \in \mathcal{S}} d_0(s) v_\pi(s) = \sum_{s \in \mathcal{S}} d_0(s) \nabla_\theta v_\pi(s)$$

then

$$\begin{aligned} \nabla_\theta \bar{v}_\pi^0 &= \sum_{s \in \mathcal{S}} d_0(s) \nabla_\theta v_\pi(s) \\ &= \sum_{s \in \mathcal{S}} d_0(s) \sum_{s' \in \mathcal{S}} \Pr_\pi(s'|s) \sum_{a \in \mathcal{A}} \nabla_\theta \pi(a|s', \theta) q_\pi(s', a) \\ &= \sum_{s' \in \mathcal{S}} \left(\sum_{s \in \mathcal{S}} d_0(s) \Pr_\pi(s'|s) \right) \sum_{a \in \mathcal{A}} \nabla_\theta \pi(a|s', \theta) q_\pi(s', a) \\ &\doteq \sum_{s' \in \mathcal{S}} \rho_\pi(s') \sum_{a \in \mathcal{A}} \nabla_\theta \pi(a|s', \theta) q_\pi(s', a) \\ &= \sum_{s \in \mathcal{S}} \rho_\pi(s) \sum_{a \in \mathcal{A}} \nabla_\theta \pi(a|s, \theta) q_\pi(s, a) \quad (\text{change } s' \text{ to } s) \\ &= \sum_{s \in \mathcal{S}} \rho_\pi(s) \sum_{a \in \mathcal{A}} \pi(a|s, \theta) \nabla_\theta \ln \pi(a|s, \theta) q_\pi(s, a) \\ &= \mathbb{E} [\nabla_\theta \ln \pi(A|S, \theta) q_\pi(S, A)], \end{aligned}$$

Theorem 1.15 (Gradient of \bar{v}_π and \bar{r}_π in the discounted case)

In the discounted case where $\gamma \in (0, 1)$, the gradient of \bar{v}_π and \bar{r}_π are

$$\nabla_\theta \bar{r}_\pi = (1 - \gamma) \nabla_\theta \bar{v}_\pi \approx \sum_{s \in \mathcal{S}} d_\pi(s) \sum_{a \in \mathcal{A}} \nabla_\theta \pi(a|s, \theta) q_\pi(s, a) = \mathbb{E} [\nabla_\theta \ln \pi(A|S, \theta) q_\pi(S, A)] \quad (1.55)$$

Here, the approximation is more accurate when γ is closer to 1.

Proof It follows from the definition of \bar{v}_π that

$$\begin{aligned} \nabla_\theta \bar{v}_\pi &= \nabla_\theta \sum_{s \in \mathcal{S}} d_\pi(s) v_\pi(s) \\ &= \sum_{s \in \mathcal{S}} \nabla_\theta d_\pi(s) v_\pi(s) + \sum_{s \in \mathcal{S}} d_\pi(s) \nabla_\theta v_\pi(s). \end{aligned}$$

This equation contains two terms. On the one hand, substituting the expression of $\nabla_\theta v_\pi$ given in (9.17) into the second term gives

$$\begin{aligned} \sum_{s \in \mathcal{S}} d_\pi(s) \nabla_\theta v_\pi(s) &= (d_\pi^\top \otimes I_m) \nabla_\theta v_\pi \\ &= (d_\pi^\top \otimes I_m) [(I_n - \gamma P_\pi)^{-1} \otimes I_m] u \\ &= [d_\pi^\top (I_n - \gamma P_\pi)^{-1}] \otimes I_m u. \end{aligned}$$

It is noted that

$$d_\pi^\top (I_n - \gamma P_\pi)^{-1} = \frac{1}{1 - \gamma} d_\pi^\top,$$

which can be easily verified by multiplying $(I_n - \gamma P_\pi)$ on both sides of the equation. Therefore, the above becomes

$$\begin{aligned} \sum_{s \in \mathcal{S}} d_\pi(s) \nabla_\theta v_\pi(s) &= \frac{1}{1 - \gamma} d_\pi^\top \otimes I_m u \\ &= \frac{1}{1 - \gamma} \sum_{s \in \mathcal{S}} d_\pi(s) \sum_{a \in \mathcal{A}} \nabla_\theta \pi(a|s, \theta) q_\pi(s, a). \end{aligned}$$

On the other hand, the first term of the earlier expression involves $\nabla_\theta d_\pi$. However, since the second term contains

$\frac{1}{1-\gamma}$, the second term becomes dominant, and the first term becomes negligible when $\gamma \rightarrow 1$. Therefore,

$$\nabla_{\theta} \bar{v}_{\pi} \approx \frac{1}{1-\gamma} \sum_{s \in \mathcal{S}} d_{\pi}(s) \sum_{a \in \mathcal{A}} \nabla_{\theta} \pi(a|s, \theta) q_{\pi}(s, a).$$

Furthermore, it follows from $\bar{r}_{\pi} = (1 - \gamma) \bar{v}_{\pi}$ that

$$\begin{aligned} \nabla_{\theta} \bar{r}_{\pi} &= (1 - \gamma) \nabla_{\theta} \bar{v}_{\pi} \approx \sum_{s \in \mathcal{S}} d_{\pi}(s) \sum_{a \in \mathcal{A}} \nabla_{\theta} \pi(a|s, \theta) q_{\pi}(s, a) \\ &= \sum_{s \in \mathcal{S}} d_{\pi}(s) \sum_{a \in \mathcal{A}} \pi(a|s, \theta) \nabla_{\theta} \ln \pi(a|s, \theta) q_{\pi}(s, a) \\ &= \mathbb{E} [\nabla_{\theta} \ln \pi(A|S, \theta) q_{\pi}(S, A)] \end{aligned}$$

For the undiscounted case, it is necessary to redefine state and action values since they maybe diverge

$$v_{\pi}(s) \doteq \mathbb{E} [(R_{t+1} - \bar{r}_{\pi}) + (R_{t+2} - \bar{r}_{\pi}) + (R_{t+3} - \bar{r}_{\pi}) + \dots | S_t = s]$$

$$q_{\pi}(s, a) \doteq \mathbb{E} [(R_{t+1} - \bar{r}_{\pi}) + (R_{t+2} - \bar{r}_{\pi}) + (R_{t+3} - \bar{r}_{\pi}) + \dots | S_t = s, A_t = a]$$

It can be verified that the state value defined above satisfies the following Bellman-like equation

$$v_{\pi}(s) = \sum_a \pi(a|s, \theta) \left[\sum_r p(r|s, a) (r - \bar{r}_{\pi}) + \sum_{s'} p(s'|s, a) v_{\pi}(s') \right] \quad (1.56)$$

The matrix form is

$$v_{\pi} = r_{\pi} - \bar{r}_{\pi} \mathbf{1}_n + P_{\pi} v_{\pi}$$

it is similar to the Bellman equation and it has a specific name called the **Poisson equation**.

Theorem 1.16 (Solution of the Poisson equation)

Let

$$v_{\pi}^* = (I_n - P_{\pi} + \mathbf{1}_n d_{\pi}^T)^{-1} r_{\pi} \quad (1.57)$$

Then, v_{π}^* is a solution of the Poisson equation. Moreover, any solution of the Poisson equation has the following form:

$$v_{\pi} = v_{\pi}^* + c \mathbf{1}_n$$

Because of c , the value of v_{π} is not unique, but \bar{r}_{π} is

$$\begin{aligned} \bar{r}_{\pi} \mathbf{1}_n &= r_{\pi} + (P_{\pi} - I_n) v_{\pi} \\ &= r_{\pi} + (P_{\pi} - I_n) (v_{\pi}^* + c \mathbf{1}_n) \\ &= r_{\pi} + (P_{\pi} - I_n) v_{\pi}^*. \end{aligned}$$

Theorem 1.17 (Gradient of \bar{r}_{π} in the undiscounted case)

$$\begin{aligned} \nabla_{\theta} \bar{r}_{\pi} &= \sum_{s \in \mathcal{S}} d_{\pi}(s) \sum_{a \in \mathcal{A}} \nabla_{\theta} \pi(a|s, \theta) q_{\pi}(s, a) \\ &= \mathbb{E} [\nabla_{\theta} \ln \pi(A|S, \theta) q_{\pi}(S, A)] \end{aligned} \quad (1.58)$$

We now introduce **Monte Carlo policy gradient** (REINFORCE), gradient-ascent algorithm for maximizing $J(\theta)$ is

$$\begin{aligned} \theta_{t+1} &= \theta_t + \alpha \nabla_{\theta} J(\theta_t) \\ &= \theta_t + \alpha \mathbb{E} \left[\nabla_{\theta} \ln \pi(A|S, \theta_t) q_{\pi}(S, A) \right] \\ &\rightarrow \text{stochastic gradient} \\ &= \theta_t + \alpha \nabla_{\theta} \ln \pi(a_t|s_t, \theta_t) q_t(s_t, a_t) \end{aligned} \quad (1.59)$$

if $q_t(s_t, a_t)$ is obtained by Monte Carlo estimation, which is one of earliest and simplest policy gradient algorithms.

rewrite it

$$\theta_{t+1} = \theta_t + \alpha \underbrace{\left(\frac{q_t(s_t, a_t)}{\pi(a_t|s_t, \theta_t)} \right)}_{\beta_t} \nabla_{\theta} \pi(a_t|s_t, \theta_t) = \theta_t + \alpha \beta_t \nabla_{\theta} \pi(a_t|s_t, \theta_t) \quad (1.60)$$

1. If $\beta_t \geq 0$, the probability of choosing (s_t, a_t) is enhanced.
2. the algorithm can strike a balance between exploration and exploitation to a certain extent due to the expression of β_t .
3. How to sample S , the distribution is long term behavior of d_{π} or ρ_{π} .
4. How to sample A , it should follow $\pi(A|S, \theta)$.
5. Unfortunately, the ideal ways for sampling S and A are not strictly followed in practice due to their low efficiency of sample usage.

Algorithm 8: Policy Gradient by Monte Carlo (REINFORCE)

Input: Initial parameter θ .

Output: Learn an optimal policy for maximizing $J(\theta)$.

```

1 foreach episode do
2   Generate an episode  $\{s_0, a_0, r_1, \dots, r_T\}$  following  $\pi(\theta)$ 
3   for  $t = 0, 1, \dots, T - 1$  do
4     Value update:  $q_t(s_t, a_t) = \sum_{k=t+1}^T \gamma^{k-t-1} r_k$ 
5     Policy update:  $\theta \leftarrow \theta + \alpha \nabla_{\theta} \ln \pi(a_t|s_t, \theta) q_t(s_t, a_t)$ 
6   end
7 end

```

1.13 Actor-Critic Methods

Here, an “actor” refers to a policy update step. The reason that it is called an actor is that the actions are taken by following the policy. Here, an “critic” refers to a value update step. It is called a critic because it criticizes the actor by evaluating its corresponding values. From another point of view, actor-critic methods are still policy gradient algorithms.

If in 1.59, we use Monte Carlo method, then it is REINFORCE, if it uses TD to evaluate $q(s, a)$, then it is **actor-critic**.

Algorithm 9: The simplest actor-critic algorithm (QAC)

Input: Initial parameter θ_0 and policy function $\pi(a|s, \theta_0)$. A value function $q(s, a, w_0)$, $\alpha_w, \alpha_{\theta} > 0$

Output: Learn an optimal policy for maximizing $J(\theta)$.

```

1 foreach time step  $t$  in each episode do
2   Generate  $a_t$  following  $\pi(a, s_t, \theta_t)$ , observe  $r_{t+1}, s_{t+1}$  and then generate  $a_{t+1}$  following  $\pi(a, s_{t+1}, \theta_t)$ 
3   Actor (Policy update):  $\theta_{t+1} \leftarrow \theta_t + \alpha_{\theta} \nabla_{\theta} \ln \pi(a_t|s_t, \theta_t) q_t(s_t, a_t, w_t)$ 
4   Critic (Value update):  $w_{t+1} = w_t + \alpha_w [r_{t+1} + \gamma q(s_{t+1}, a_{t+1}, w_t) - q(s_t, a_t, w_t)] \nabla_w q(s_t, a_t, w_t)$ 
5 end

```

We then introduce **Advantage Actor-Critic** (A2C), the core idea of this algorithm is to introduce a baseline to reduce estimation variance.

One interesting property of the policy gradient is that it is invariant to an additional baseline

$$\mathbb{E}_{S \sim \eta, A \sim \pi} [\nabla_{\theta} \ln \pi(A|S, \theta_t) q_{\pi}(S, A)] = \mathbb{E}_{S \sim \eta, A \sim \pi} [\nabla_{\theta} \ln \pi(A|S, \theta_t) (q_{\pi}(S, A) - b(S))] \quad (1.61)$$

this holds if and only if

$$\mathbb{E}_{S \sim \eta, A \sim \pi} [\nabla_{\theta} \ln \pi(A|S, \theta_t) b(S)] = 0$$

this is valid because

$$\begin{aligned}
\mathbb{E}_{S \sim \eta, A \sim \pi} [\nabla_{\theta} \ln \pi(A|S, \theta_t) b(S)] &= \sum_{s \in \mathcal{S}} \eta(s) \sum_{a \in \mathcal{A}} \pi(a|s, \theta_t) \nabla_{\theta} \ln \pi(a|s, \theta_t) b(s) \\
&= \sum_{s \in \mathcal{S}} \eta(s) \sum_{a \in \mathcal{A}} \nabla_{\theta} \pi(a|s, \theta_t) b(s) \\
&= \sum_{s \in \mathcal{S}} \eta(s) b(s) \sum_{a \in \mathcal{A}} \nabla_{\theta} \pi(a|s, \theta_t) \\
&= \sum_{s \in \mathcal{S}} \eta(s) b(s) \nabla_{\theta} \sum_{a \in \mathcal{A}} \pi(a|s, \theta_t) \\
&= \sum_{s \in \mathcal{S}} \eta(s) b(s) \nabla_{\theta} 1 = 0
\end{aligned}$$

The baseline is useful because it can reduce the approximation variance when we use samples to approximate the true gradient. In particular, let

$$X(S, A) \doteq \nabla_{\theta} \ln \pi(A|S, \theta_t) [q_{\pi}(S, A) - b(S)] \quad (1.62)$$

Then, the true gradient is $\mathbb{E}[X(S, A)]$. In REINFORCE and QAC, we set $b = 0$, which is not guaranteed to be a good baseline.

In fact, the optimal baseline that minimizes $\text{var}(X)$ is (please refer to book box 10.1)

$$b^*(s) = \frac{\mathbb{E}_{A \sim \pi} [||\nabla_{\theta} \ln \pi(A|s, \theta_t)||^2 q_{\pi}(s, A)]}{\mathbb{E}_{A \sim \pi} [||\nabla_{\theta} \ln \pi(A|s, \theta_t)||^2]}, \quad s \in \mathcal{S} \quad (1.63)$$

but this is too complex, we use a suboptimal one:

$$b^{\dagger}(s) = \mathbb{E}_{A \sim \pi} [q_{\pi}(s, A)] = v_{\pi}(s), \quad s \in \mathcal{S}$$

Interestingly, this suboptimal baseline is the state value.

The algorithm is then

$$\begin{aligned}
\theta_{t+1} &= \theta_t + \alpha \mathbb{E} [\nabla_{\theta} \ln \pi(A|S, \theta_t) [q_{\pi}(S, A) - v_{\pi}(S)]] \\
&= \theta_t + \alpha \mathbb{E} [\nabla_{\theta} \ln \pi(A|S, \theta_t) \delta_{\pi}(S, A)] \\
&\rightarrow \text{stochastic version} \\
&= \theta_t + \alpha \nabla_{\theta} \ln \pi(a_t|s_t, \theta_t) [q_t(s_t, a_t) - v_t(s_t)] \\
&= \theta_t + \alpha \nabla_{\theta} \ln \pi(a_t|s_t, \theta_t) \delta_t(s_t, a_t)
\end{aligned} \quad (1.64)$$

If $q_t(s_t, a_t)$ and $v_t(s_t)$ are estimated by Monte Carlo learning, this is called REINFORCE with a baseline. If they are estimated by TD learning, the algorithm is usually called advantage actor-critic (A2C).

Algorithm 10: Advantage actor-critic (A2C) or TD actor-critic

Input: Initial parameter θ_0 and policy function $\pi(a|s, \theta_0)$. A value function $v(s, w_0)$, $\alpha_w, \alpha_{\theta} > 0$

Output: Learn an optimal policy for maximizing $J(\theta)$.

```

1 foreach time step  $t$  in each episode do
2   Generate  $a_t$  following  $\pi(a, s_t, \theta_t)$ , observe  $r_{t+1}, s_{t+1}$ 
3   Advantage:  $\delta_t = r_{t+1} + \gamma v_t(s_{t+1}, w_t) - v_t(s_t, w_t)$ 
4   Actor (Policy update):  $\theta_{t+1} \leftarrow \theta_t + \alpha_{\theta} \delta_t \nabla_{\theta} \ln \pi(a_t|s_t, \theta_t)$ 
5   Critic (Value update):  $w_{t+1} = w_t + \alpha_w \delta_t \nabla_w v(s_t, w_t)$ 
6 end

```

It should be noted that the advantage function in this implementation is approximated by the TD error

$$q_t(s_t, a_t) - v_t(s_t) \approx r_{t+1} + \gamma v_t(s_{t+1}) - v_t(s_t)$$

This approximation is reasonable because

$$q_{\pi}(s_t, a_t) - v_{\pi}(s_t) = \mathbb{E} [R_{t+1} + \gamma v_{\pi}(S_{t+1}) - v_{\pi}(S_t) | S_t = s_t, A_t = a_t]$$

When we use the TD error, the algorithm may also be called TD actor-critic.

All the methods we talked above is on-policy, now we introduce off-policy one. But first we introduce *Importance Sampling*, if x is generated with p_1 , we want use p_0 to evaluate expectation

$$\mathbb{E}_{X \sim p_0}[X] = \sum_{x \in \mathcal{X}} p_0(x)x = \sum_{x \in \mathcal{X}} p_1(x) \underbrace{\frac{p_0(x)}{p_1(x)}}_{f(x)x} = \mathbb{E}_{X \sim p_1}[f(X)] \quad (1.65)$$

and then

$$\mathbb{E}_{X \sim p_0}[X] = \mathbb{E}_{X \sim p_1}[f(X)] \approx \bar{f} = \frac{1}{n} \sum_{i=1}^n f(x_i) = \frac{1}{n} \sum_{i=1}^n \underbrace{\frac{p_0(x_i)}{p_1(x_i)}}_{\text{importance weight}} x_i. \quad (1.66)$$

Suppose the behavior policy is β , than what we want to optimize is

$$J(\theta) = \sum_{s \in \mathcal{S}} d_\beta(s) v_\pi(s) = \mathbb{E}_{S \sim d_\beta}[v_\pi(s)]$$

Theorem 1.18 (Off-policy policy gradient theorem)

In the discounted case, the gradient is

$$\nabla_\theta J(\theta) = \mathbb{E}_{S \sim \rho, A \sim \beta} \left[\underbrace{\frac{\pi(A|S, \theta)}{\beta(A|S)}}_{\text{importance weight}} \nabla_\theta \ln \pi(A|S, \theta) q_\pi(S, A) \right] \quad (1.67)$$

where the state distribution ρ is

$$\rho(s) \doteq \sum_{s' \in \mathcal{S}} d_\beta(s') \Pr_\pi(s|s'), \quad s \in \mathcal{S}$$

where

$$\Pr_\pi(s|s') = \sum_{k=0}^{\infty} \gamma^k [P_\pi^k]_{s's} = [(I - \gamma P_\pi)^{-1}]_{s's}$$

Proof Since d_β is independent of θ , the gradient of $J(\theta)$ satisfies

$$\nabla_\theta J(\theta) = \nabla_\theta \sum_{s \in \mathcal{S}} d_\beta(s) v_\pi(s) = \sum_{s \in \mathcal{S}} d_\beta(s) \nabla_\theta v_\pi(s). \quad (1.68)$$

According to Lemma 1.4, the expression of $\nabla_\theta v_\pi(s)$ is

$$\nabla_\theta v_\pi(s) = \sum_{s' \in \mathcal{S}} \Pr_\pi(s'|s) \sum_{a \in \mathcal{A}} \nabla_\theta \pi(a|s', \theta) q_\pi(s', a) \quad (1.69)$$

where $\Pr_\pi(s'|s) \doteq \sum_{k=0}^{\infty} \gamma^k [P_\pi^k]_{ss'} = [(I_n - \gamma P_\pi)^{-1}]_{ss'}$. Then we have

$$\begin{aligned} \nabla_\theta J(\theta) &= \sum_{s \in \mathcal{S}} d_\beta(s) \nabla_\theta v_\pi(s) = \sum_{s \in \mathcal{S}} d_\beta(s) \sum_{s' \in \mathcal{S}} \Pr_\pi(s'|s) \sum_{a \in \mathcal{A}} \nabla_\theta \pi(a|s', \theta) q_\pi(s', a) \\ &= \sum_{s' \in \mathcal{S}} \left(\sum_{s \in \mathcal{S}} d_\beta(s) \Pr_\pi(s'|s) \right) \sum_{a \in \mathcal{A}} \nabla_\theta \pi(a|s', \theta) q_\pi(s', a) \\ &\doteq \sum_{s' \in \mathcal{S}} \rho(s') \sum_{a \in \mathcal{A}} \nabla_\theta \pi(a|s', \theta) q_\pi(s', a) \\ &= \sum_{s \in \mathcal{S}} \rho(s) \sum_{a \in \mathcal{A}} \nabla_\theta \pi(a|s, \theta) q_\pi(s, a) \quad (\text{change } s' \text{ to } s) \\ &= \mathbb{E}_{S \sim \rho} \left[\sum_{a \in \mathcal{A}} \nabla_\theta \pi(a|S, \theta) q_\pi(S, a) \right]. \end{aligned}$$

By using the importance sampling technique, the above equation can be further rewritten as

$$\begin{aligned}\mathbb{E}_{S \sim \rho} \left[\sum_{a \in \mathcal{A}} \nabla_{\theta} \pi(a|S, \theta) q_{\pi}(S, a) \right] &= \mathbb{E}_{S \sim \rho} \left[\sum_{a \in \mathcal{A}} \beta(a|S) \frac{\pi(a|S, \theta)}{\beta(a|S)} \frac{\nabla_{\theta} \pi(a|S, \theta)}{\pi(a|S, \theta)} q_{\pi}(S, a) \right] \\ &= \mathbb{E}_{S \sim \rho} \left[\sum_{a \in \mathcal{A}} \beta(a|S) \frac{\pi(a|S, \theta)}{\beta(a|S)} \nabla_{\theta} \ln \pi(a|S, \theta) q_{\pi}(S, a) \right] \\ &= \mathbb{E}_{S \sim \rho, A \sim \beta} \left[\frac{\pi(A|S, \theta)}{\beta(A|S)} \nabla_{\theta} \ln \pi(A|S, \theta) q_{\pi}(S, A) \right].\end{aligned}$$

The corresponding algorithm is

$$\theta_{t+1} = \theta_t + \alpha_{\theta} \frac{\pi(a_t|s_t, \theta_t)}{\beta(a_t|s_t)} \nabla_{\theta} \ln \pi(a_t|s_t, \theta_t) (q_t(s_t, a_t) - v_t(s_t))$$

using TD error, then we get

$$\theta_{t+1} = \theta_t + \alpha_{\theta} \frac{\pi(a_t|s_t, \theta)}{\beta(a_t|s_t)} \nabla_{\theta} \ln \pi(a_t|s_t, \theta) \delta_t(s_t, a_t)$$

where

$$q_t(s_t, a_t) - v_t(s_t) \approx r_{t+1} + \gamma v_t(s_{t+1}) - v_t(s_t) \doteq \delta_t(s_t, a_t)$$

Algorithm 11: Off-policy actor-critic based on importance sampling

Input: A given behavior policy $\beta(a|s)$. Initial parameter θ_0 and policy function $\pi(a|s, \theta_0)$. A value function $v(s, w_0)$, $\alpha_w, \alpha_{\theta} > 0$

Output: Learn an optimal policy for maximizing $J(\theta)$.

```

1 foreach time step  $t$  in each episode do
2   Generate  $a_t$  following  $\beta(s_t)$ , observe  $r_{t+1}, s_{t+1}$ 
3   Advantage:  $\delta_t = r_{t+1} + \gamma v_t(s_{t+1}, w_t) - v_t(s_t, w_t)$ 
4   Actor (Policy update):  $\theta_{t+1} = \theta_t + \alpha_{\theta} \frac{\pi(a_t|s_t, \theta_t)}{\beta(a_t|s_t)} \delta_t \nabla_{\theta} \ln \pi(a_t|s_t, \theta_t)$ 
5   Critic (Value update):  $w_{t+1} = w_t + \alpha_w \frac{\pi(a_t|s_t, \theta_t)}{\beta(a_t|s_t)} \delta_t \nabla_w v(s_t, w_t)$ 
6 end
```

Up to now, we need $\pi(a|s, \theta) > 0$ because the use of \ln , we will then show deterministic policy can also be used in policy gradient methods. We use $a = \mu(s, \theta)$ to specifically denote a deterministic policy.

Theorem 1.19 (Deterministic policy gradient theorem)

The gradient of $J(\theta)$ is

$$\begin{aligned}\nabla_{\theta} J(\theta) &= \sum_{s \in \mathcal{S}} \eta(s) \nabla_{\theta} \mu(s) (\nabla_a q_{\mu}(s, a)) \Big|_{a=\mu(s)} \\ &= \mathbb{E}_{S \sim \eta} [\nabla_{\theta} \mu(S) (\nabla_a q_{\mu}(S, a)) \Big|_{a=\mu(S)}]\end{aligned}\tag{1.70}$$

where η is a distribution of the states.

When we use samples to approximate the true gradient, it is not required to sample actions. Therefore, the deterministic policy gradient method is off-policy.

We first derive the gradient of the average value:

$$J(\theta) = \mathbb{E}[v_{\mu}(s)] = \sum_{s \in \mathcal{S}} d_0(s) v_{\mu}(s)\tag{1.71}$$

where d_0 is the probability distribution of the states. Here, d_0 is selected to be independent of μ for simplicity.

Lemma 1.5 (Gradient of $v_\mu(s)$)

In the discounted case, it holds for any $s \in \mathcal{S}$ that

$$\nabla_\theta v_\mu(s) = \sum_{s' \in \mathcal{S}} \Pr_\mu(s'|s) \nabla_\theta \mu(s') (\nabla_a q_\mu(s', a))|_{a=\mu(s')} \quad (1.72)$$

where

$$\Pr_\mu(s'|s) \doteq \sum_{k=0}^{\infty} \gamma^k [P_\mu^k]_{ss'} = [(I - \gamma P_\mu)^{-1}]_{ss'}$$

is the discounted total probability of transitioning from s to s' under μ .

**Theorem 1.20 (Deterministic policy gradient theorem in the discounted case)**

The distribution η in 1.19 is

$$\rho_\mu(s) = \sum_{s' \in \mathcal{S}} d_0(s') \Pr_\mu(s|s')$$

where

$$\Pr_\mu(s'|s) \doteq \sum_{k=0}^{\infty} \gamma^k [P_\mu^k]_{ss'} = [(I - \gamma P_\mu)^{-1}]_{ss'}$$

is the discounted total probability of transitioning from s to s' under μ .



We next derive the gradient of the average reward

$$J(\theta) = \bar{r}_\mu = \mathbb{E}_{S \sim d_\mu} [r_\mu(S)] = \sum_{s \in \mathcal{S}} d_\mu(s) r_\mu(s) \quad (1.73)$$

where

$$r_\mu(s) = \mathbb{E}[R|s, a = \mu(s)] = \sum_r r p(r|s, a = \mu(s))$$

Theorem 1.21 (Deterministic policy gradient theorem in the undiscounted case)

The gradient of $J(\theta)$ is

$$\begin{aligned} \nabla_\theta J(\theta) &= \sum_{s \in \mathcal{S}} d_\mu(s) \nabla_\theta \mu(s) (\nabla_a q_\mu(s, a))|_{a=\mu(s)} \\ &= \mathbb{E}_{S \sim d_\mu} [\nabla_\theta \mu(S) (\nabla_a q_\mu(S, a))|_{a=\mu(S)}] \end{aligned} \quad (1.74)$$

where d_μ is a stationary distribution of the states.



So we can design a gradient ascent algorithm

$$\theta_{t+1} = \theta_t + \alpha_\theta \mathbb{E}_{S \sim \eta} [\nabla_\theta \mu(S) (\nabla_a q_\mu(S, a))|_{a=\mu(S)}]$$

the stochastic version is

$$\theta_{t+1} = \theta_t + \alpha_\theta \nabla_\theta \mu(s_t) (\nabla_a q_\mu(s_t, a))|_{a=\mu(s_t)}$$

In particular, the experience sample required by the critic is $(s_t, \alpha_t, r_{t+1}, s_{t+1}, \tilde{a}_{t+1})$, where $\tilde{a}_{t+1} = \mu(s_{t+1})$. The generation of this experience sample involves two policies. The first is the policy for generating a_t at s_t , and the second is the policy for generating \tilde{a}_{t+1} at s_{t+1} . The first policy that generates a_t is the behavior policy since a_t is used to interact with the environment. The second policy must be μ because it is the policy that the critic aims to evaluate. Hence, μ is the target policy. It should be noted that \tilde{a}_{t+1} is not used to interact with the environment in the next time step. Hence, μ is not the behavior policy. Therefore, the critic is off-policy. And we showed above the actor part is also off-policy.

Algorithm 12: Deterministic policy gradient or deterministic actor-critic

Input: A given behavior policy $\beta(a|s)$. Initial parameter θ_0 and policy function $\mu(s, \theta_0)$. A value function $q(s, a, w_0)$, $\alpha_w, \alpha_\theta > 0$

Output: Learn an optimal policy for maximizing $J(\theta)$.

```

1 foreach time step  $t$  in each episode do
2   Generate  $a_t$  following  $\beta(s_t)$ , observe  $r_{t+1}, s_{t+1}$ 
3   Advantage:  $\delta_t = r_{t+1} + \gamma q(s_{t+1}, \mu(s_{t+1}, \theta_t), w_t) - q(s_t, a_t, w_t)$ 
4   Actor (Policy update):  $\theta_{t+1} = \theta_t + \alpha_\theta \nabla_\theta \mu(s_t, \theta_t) (\nabla_a q_\mu(s_t, a, w_t))|_{a=\mu(s_t)}$ 
5   Critic (Value update):  $w_{t+1} = w_t + \alpha_w \delta_t \nabla_w q(s_t, a_t, w_t)$ 
6 end

```

1.14 Appendix: Measure-Theoretic Probability Theory

A **Probability Triple** is a tuple (Ω, \mathcal{F}, P) , where:

- Ω is a sample space, which is a set of all possible outcomes;
- \mathcal{F} is a σ -algebra, which is a collection of subsets of Ω that includes the empty set and Ω , and is closed under complementation and countable unions; An element in \mathcal{F} is called an event, denoted as A . An elementary event refers to a single outcome in Ω .
- P is a probability measure, which assigns a probability to each event in \mathcal{F} such that $P(\Omega) = 1$ and $P(\emptyset) = 0$.

For example, the game of dice, $\Omega = \{1, 2, 3, 4, 5, 6\}$, $\mathcal{F} = \mathcal{P}(\Omega)$ is the power set of Ω , we can define one event $A = \{w \in \Omega : w > 3\} = \{4, 5, 6\}$ with $\mathbb{P}(A) = 1/2$.

Definition 1.7 (σ -algebra)

A **σ -algebra** \mathcal{F} on a set Ω is a collection of subsets of Ω that satisfies the following properties:

1. $\emptyset \in \mathcal{F}$ and $\Omega \in \mathcal{F}$;
2. If $A \in \mathcal{F}$, then $A^c = \Omega \setminus A \in \mathcal{F}$ (closed under complementation);
3. If $A_1, A_2, A_3, \dots \in \mathcal{F}$, then $\bigcup_{i=1}^{\infty} A_i \in \mathcal{F}$ (closed under countable unions).



A **Random Variable** is a measurable function $X(w) : \Omega \rightarrow \mathbb{R}$ that maps **each** outcome in Ω to real numbers. Not all the mappings are random variables, the formal definition is:

$$A = \{w \in \Omega \mid X(w) \leq x\} \in \mathcal{F}, \forall x \in \mathbb{R}$$

This definition indicates that X is a random variable only if $X(w) \leq x$ is an event in \mathcal{F} , it is probabilistically measurable.

Simple random variables are those that take a finite number of values, and they can be expressed as:

$$X(w) = \sum_{x \in \mathcal{X}} x \mathbb{I}_{A_x}(w)$$

where $A_x = \{w \in \Omega \mid X(w) = x\} = X^{-1}(x)$ and

$$\mathbb{I}_{A_x}(w) = \begin{cases} 1 & \text{if } w \in A_x, \\ 0 & \text{otherwise.} \end{cases}$$

With the above preparation, we can define the **Expectation** of a random variable X :

$$\mathbb{E}[X] = \sum_{x \in \mathcal{X}} x \mathbb{P}(A_x) \quad (1.75)$$

Lemma 1.6 (Basic Properties)

1. $\mathbb{E}[a \mid Y] = a$, where a is a constant;
2. $\mathbb{E}[aX + bZ \mid Y] = a\mathbb{E}[X \mid Y] + b\mathbb{E}[Z \mid Y]$, where a, b are constants;

3. $\mathbb{E}[X | Y] = \mathbb{E}[X]$ if X is independent of Y ;
4. $\mathbb{E}[X f(Y) | Y] = f(Y) \mathbb{E}[X | Y]$, where $f(Y)$ is a function of Y ;
5. $\mathbb{E}[f(Y) | Y] = f(Y)$;
6. $\mathbb{E}[X | Y, f(Y)] = \mathbb{E}[X | Y]$;
7. If $X \geq 0$, then $\mathbb{E}[X | Y] \geq 0$;
8. If $X \geq Z$, then $\mathbb{E}[X | Y] \geq \mathbb{E}[Z | Y]$;
9. $\mathbb{E}[\mathbb{E}[X | Y]] = \mathbb{E}[X]$;
10. $\mathbb{E}[\mathbb{E}[X | Y, Z]] = \mathbb{E}[X]$;
11. $\mathbb{E}[\mathbb{E}[X | Y] | Y] = \mathbb{E}[X | Y]$;



Given a stochastic sequence $\{X_k\} = \{X_1, X_2, \dots\}$:

Definition 1.8 (Sure Convergence)

$\{X_k\}$ converges **surely** to X if:

$$\lim_{k \rightarrow \infty} X_k(w) = X(w), \quad \forall w \in \Omega$$

this can be equivalently written as:

$$A = \Omega \quad \text{where} \quad A = \{w \in \Omega \mid \lim_{k \rightarrow \infty} X_k(w) = X(w)\}$$



Definition 1.9 (Almost Sure Convergence)

$\{X_k\}$ converges **almost surely** (or almost everywhere or with probability 1) to X if:

$$\mathbb{P}(A) = 1 \quad \text{where} \quad A = \{w \in \Omega \mid \lim_{k \rightarrow \infty} X_k(w) = X(w)\}$$

The points, for which this limit is invalid, form a set of zero measure. For the sake of simplicity, it is often written as:

$$P\left(\lim_{k \rightarrow \infty} X_k = X\right) = 1$$

and denote as $X_k \xrightarrow{a.s.} X$.



Definition 1.10 (Convergence in Probability)

$\{X_k\}$ converges **in probability** to X if $\forall \epsilon > 0$:

$$\lim_{k \rightarrow \infty} \mathbb{P}(A_k) = 0 \quad \text{where} \quad A_k = \{w \in \Omega : |X_k(w) - X(w)| > \epsilon\}$$

This means that the probability of the difference between X_k and X being greater than ϵ goes to zero as k increases. It can be written as:

$$\lim_{k \rightarrow \infty} \mathbb{P}(|X_k - X| > \epsilon) = 0$$



Definition 1.11 (Convergence in Mean)

$\{X_k\}$ converges **in mean** (r -th mean or in L^r norm) to X if:

$$\lim_{k \rightarrow \infty} \mathbb{E}[|X_k - X|^r] = 0$$

This means that the expected value of the absolute difference between X_k and X goes to zero as k increases.



Definition 1.12 (Convergence in Distribution)

$\{X_k\}$ converges **in distribution** to X if:

$$\lim_{k \rightarrow \infty} \mathbb{P}(X_k \leq a) = \mathbb{P}(X \leq a) \quad \forall a \in \mathbb{R}$$

a compact expression is:

$$\lim_{k \rightarrow \infty} \mathbb{P}(A_k) = \mathbb{P}(A)$$

where

$$A_k = \{w \in \Omega : X_k(w) \leq a\}, A = \{w \in \Omega : X(w) \leq a\}$$



The relationship between these convergence types is: almost sure convergence \rightarrow convergence in probability \rightarrow convergence in distribution. Convergence in mean \rightarrow convergence in probability \rightarrow convergence in distribution. And convergence in mean does not imply almost sure convergence.

Chapter 2 From LQR to RL

Introduction

□ LQR Problem

□ Reinforcement Learning

□ iLQR and DDP

2.1 LQR and Value function

Given a linear model $x_{t+1} = f(x_t, u_k) = A_t x_t + B_t u_t + C_t$. We want to optimize:

$$\min_{u_1, \dots, u_T} c(x_1, u_1) + c(f(x_1, u_1), u_2) + \dots + c(f(f(\dots)), u_T)$$

where we denotes

$$c(x_t, u_t) = \frac{1}{2} \begin{bmatrix} x_t \\ u_t \end{bmatrix}^T C_t \begin{bmatrix} x_t \\ u_t \end{bmatrix} + \begin{bmatrix} x_t \\ u_t \end{bmatrix}^T c_t$$

and

$$f(x_t, u_t) = F_t \begin{bmatrix} x_t \\ u_t \end{bmatrix} + f_t$$

We first do the **Backward Recursion**, solve for u_T only, then the action value function (or the negative cost function, here we take them with same sign) is:

$$Q(x_T, u_T) = \text{const} + \frac{1}{2} \begin{bmatrix} x_T \\ u_T \end{bmatrix}^T C_T \begin{bmatrix} x_T \\ u_T \end{bmatrix} + \begin{bmatrix} x_T \\ u_T \end{bmatrix}^T c_T$$

Get the derivative respect to u_T , which is:

$$\nabla_{u_T} Q(x_T, u_T) = C_{u_T, x_T} x_T + C_{u_T, u_T} u_T + c_{u_T}^T = 0$$

so we can get $K_T = -C_{u_T, u_T}^{-1} C_{u_T, x_T}$, $k_T = -C_{u_T, u_T}^{-1} c_{u_T}$.

And we get the policy (which is a linear policy): $u_T = K_T x_T + k_T$. Because u_T is fully determined by x_T , we can eliminate it via substitution:

$$V(x_T) = \text{const} + \frac{1}{2} \begin{bmatrix} x_T \\ K_T x_T + k_T \end{bmatrix}^T C_T \begin{bmatrix} x_T \\ K_T x_T + k_T \end{bmatrix} + \begin{bmatrix} x_T \\ K_T x_T + k_T \end{bmatrix}^T c_T$$

Open the equation:

$$\begin{aligned} V(x_T) &= \text{const} + \frac{1}{2} \mathbf{x}_T^T \mathbf{V}_T \mathbf{x}_T + \mathbf{x}_T^T \mathbf{v}_T \\ \mathbf{V}_T &= \mathbf{C}_{\mathbf{x}_T, \mathbf{x}_T} + \mathbf{C}_{\mathbf{x}_T, \mathbf{u}_T} \mathbf{K}_T + \mathbf{K}_T^T \mathbf{C}_{\mathbf{u}_T, \mathbf{x}_T} + \mathbf{K}_T^T \mathbf{C}_{\mathbf{u}_T, \mathbf{u}_T} \mathbf{K}_T \\ \mathbf{v}_T &= \mathbf{c}_{\mathbf{x}_T} + \mathbf{C}_{\mathbf{x}_T, \mathbf{u}_T} \mathbf{k}_T + \mathbf{K}_T^T \mathbf{c}_{\mathbf{u}_T} + \mathbf{K}_T^T \mathbf{C}_{\mathbf{u}_T, \mathbf{u}_T} \mathbf{k}_T \end{aligned}$$

Use x_{T-1} and u_{T-1} to substitute the action value equation:

$$Q(x_{T-1}, u_{T-1}) = \text{const} + \frac{1}{2} \begin{bmatrix} x_{T-1} \\ u_{T-1} \end{bmatrix}^T C_{T-1} \begin{bmatrix} x_{T-1} \\ u_{T-1} \end{bmatrix} + \begin{bmatrix} x_{T-1} \\ u_{T-1} \end{bmatrix}^T c_{T-1} + V(f(x_{T-1}, u_{T-1}))$$

And the value function of x_T can be written as:

$$V(x_T) = \text{const} + \frac{1}{2} \begin{bmatrix} x_{T-1} \\ u_{T-1} \end{bmatrix}^T F_{T-1}^T V_T F_{T-1} \begin{bmatrix} x_{T-1} \\ u_{T-1} \end{bmatrix} + \begin{bmatrix} x_{T-1} \\ u_{T-1} \end{bmatrix}^T F_{T-1}^T V_T f_{T-1} + \begin{bmatrix} x_{T-1} \\ u_{T-1} \end{bmatrix}^T F_{T-1}^T \mathbf{v}_T$$

So the action value function is:

$$Q(x_{T-1}, u_{T-1}) = \text{const} + \frac{1}{2} \begin{bmatrix} x_{T-1} \\ u_{T-1} \end{bmatrix}^T Q_{T-1} \begin{bmatrix} x_{T-1} \\ u_{T-1} \end{bmatrix} + \begin{bmatrix} x_{T-1} \\ u_{T-1} \end{bmatrix}^T q_{T-1}$$

where

$$\begin{aligned} Q_{T-1} &= C_{T-1} + F_{T-1}^T V_T F_{T-1} \\ q_{T-1} &= c_{T-1} + F_{T-1}^T V_T f_{T-1} + F_{T-1}^T v_T \end{aligned}$$

get the derivative:

$$\nabla_{u_{T-1}} Q(x_{T-1}, u_{T-1}) = Q_{u_{T-1}, x_{T-1}} x_{T-1} + Q_{u_{T-1}, u_{T-1}} u_{T-1} + q_{u_{T-1}}^T = 0$$

where

$$\begin{aligned} u_{T-1} &= K_{T-1} x_{T-1} + k_{T-1} \\ K_{T-1} &= -Q_{u_{T-1}, u_{T-1}}^{-1} Q_{u_{T-1}, x_{T-1}} \\ k_{T-1} &= -Q_{u_{T-1}, u_{T-1}}^{-1} q_{u_{T-1}} \end{aligned}$$

So we can continue the substitution process to the first control u_0 , with the information of x_0 , we start a **Forward Recursion**, $u_t = K_t x_t + k_t$, $x_{t+1} = f(x_t, u_t)$ to get the future states.

We can write the Backward Recursion as algorithm:

Algorithm 13: Backward Pass for Value Function Computation

```

1 for  $t = T$  to 1 do
2    $\mathbf{Q}_t = \mathbf{C}_t + \mathbf{F}_t^T \mathbf{V}_{t+1} \mathbf{F}_t$ ;
3    $\mathbf{q}_t = \mathbf{c}_t + \mathbf{F}_t^T \mathbf{V}_{t+1} \mathbf{f}_t + \mathbf{F}_t^T \mathbf{v}_{t+1}$ ;
4    $Q(x_t, u_t) = \text{const} + \frac{1}{2} \begin{bmatrix} x_t \\ u_t \end{bmatrix}^T \mathbf{Q}_t \begin{bmatrix} x_t \\ u_t \end{bmatrix} + \begin{bmatrix} x_t \\ u_t \end{bmatrix}^T \mathbf{q}_t$ ;
5    $u_t \leftarrow \arg \min_{u_t} Q(x_t, u_t) = \mathbf{K}_t x_t + \mathbf{k}_t$ ;
6    $\mathbf{K}_t = -\mathbf{Q}_{u_t, u_t}^{-1} \mathbf{Q}_{u_t, x_t}$ ;
7    $\mathbf{k}_t = -\mathbf{Q}_{u_t, u_t}^{-1} \mathbf{q}_{u_t}$ ;
8    $\mathbf{V}_t = \mathbf{Q}_{x_t, x_t} + \mathbf{Q}_{x_t, u_t} \mathbf{K}_t + \mathbf{K}_t^T \mathbf{Q}_{u_t, x_t} + \mathbf{K}_t^T \mathbf{Q}_{u_t, u_t} \mathbf{K}_t$ ;
9    $\mathbf{v}_t = \mathbf{q}_{x_t} + \mathbf{Q}_{x_t, u_t} \mathbf{k}_t + \mathbf{K}_t^T \mathbf{q}_{u_t} + \mathbf{K}_t^T \mathbf{Q}_{u_t, u_t} \mathbf{k}_t$ ;
10   $V(x_t) = \text{const} + \frac{1}{2} x_t^T \mathbf{V}_t x_t + x_t^T \mathbf{v}_t$ ;
11 end
```

We can generalize it to stochastic case, where system dynamic with a gaussian noise (control is still deterministic), because the expectation of gaussian is zero for linear and constant for quadratic cost ($\mathbb{E}[x_{t+1}^T V x_{t+1}] = (Ax_t + Bu_t)^T V (Ax_t + Bu_t) + \text{tr}(VW)$, so when minimizing it is ignored).

For the nonlinear case, we linearize it with reference point:

$$\begin{aligned} f(\mathbf{x}_t, \mathbf{u}_t) &\approx f(\hat{\mathbf{x}}_t, \hat{\mathbf{u}}_t) + \nabla_{\mathbf{x}_t, \mathbf{u}_t} f(\hat{\mathbf{x}}_t, \hat{\mathbf{u}}_t) \begin{bmatrix} \mathbf{x}_t - \hat{\mathbf{x}}_t \\ \mathbf{u}_t - \hat{\mathbf{u}}_t \end{bmatrix} \\ c(\mathbf{x}_t, \mathbf{u}_t) &\approx c(\hat{\mathbf{x}}_t, \hat{\mathbf{u}}_t) + \nabla_{\mathbf{x}_t, \mathbf{u}_t} c(\hat{\mathbf{x}}_t, \hat{\mathbf{u}}_t) \begin{bmatrix} \mathbf{x}_t - \hat{\mathbf{x}}_t \\ \mathbf{u}_t - \hat{\mathbf{u}}_t \end{bmatrix} \\ &\quad + \frac{1}{2} \begin{bmatrix} \mathbf{x}_t - \hat{\mathbf{x}}_t \\ \mathbf{u}_t - \hat{\mathbf{u}}_t \end{bmatrix}^T \nabla_{\mathbf{x}_t, \mathbf{u}_t}^2 c(\hat{\mathbf{x}}_t, \hat{\mathbf{u}}_t) \begin{bmatrix} \mathbf{x}_t - \hat{\mathbf{x}}_t \\ \mathbf{u}_t - \hat{\mathbf{u}}_t \end{bmatrix} \end{aligned}$$

and use

$$\bar{f}(\delta \mathbf{x}_t, \delta \mathbf{u}_t) = \mathbf{F}_t \begin{bmatrix} \delta \mathbf{x}_t \\ \delta \mathbf{u}_t \end{bmatrix}, \quad \bar{c}(\delta \mathbf{x}_t, \delta \mathbf{u}_t) = \frac{1}{2} \begin{bmatrix} \delta \mathbf{x}_t \\ \delta \mathbf{u}_t \end{bmatrix}^T \mathbf{C}_t \begin{bmatrix} \delta \mathbf{x}_t \\ \delta \mathbf{u}_t \end{bmatrix} + \begin{bmatrix} \delta \mathbf{x}_t \\ \delta \mathbf{u}_t \end{bmatrix}^T \mathbf{c}_t$$

Now we can run LQR with it, this is called **iLQR**, we get $u_t = K_t(x_t - \hat{x}_t) + k_t + \hat{u}_t$, this is an approximation of Newton's method for solving the entire cost function over the horizon. And if we linearize the dynamic with second order information:

$$f(\mathbf{x}_t, \mathbf{u}_t) \approx f(\hat{\mathbf{x}}_t, \hat{\mathbf{u}}_t) + \nabla_{\mathbf{x}_t, \mathbf{u}_t} f(\hat{\mathbf{x}}_t, \hat{\mathbf{u}}_t) \begin{bmatrix} \delta \mathbf{x}_t \\ \delta \mathbf{u}_t \end{bmatrix} + \frac{1}{2} \left(\nabla_{\mathbf{x}_t, \mathbf{u}_t}^2 f(\hat{\mathbf{x}}_t, \hat{\mathbf{u}}_t) \begin{bmatrix} \delta \mathbf{x}_t \\ \delta \mathbf{u}_t \end{bmatrix} \right) \begin{bmatrix} \delta \mathbf{x}_t \\ \delta \mathbf{u}_t \end{bmatrix}$$

it is called *DDP*.