

## SD = REINFORCEMENT LEARNING

"Reinforcement learning is learning what to do - how to map situations to actions - so as to maximize a numerical reward signal"  $\equiv$  "learning from interactions to achieve goals"

[R.S. Sutton & A.G. Barto. Reinforcement Learning = An Introduction.]

↳ In Reinforcement Learning (RL) problem, there is no supervisor, only a reward ( $\neq$  from Supervised Learning) + data is sequential: time really matters.

The learning agent's actions influence the subsequent data it receives  $\rightarrow$  he should be able to sense the state of his environment, and must be able to take actions that affect the state. The agent must also have goals related to the state of the environment. The mathematical model that captures these aspects is known as Markov Decision Process (MDP).

Ex: (i) Defeat the world champion at Backgammon / Go  
(ii) Make a robot walk  
(iii) Manage an investment portfolio.  
(iv) Adaptive Controller adjusts parameters of a petroleum refinery's operation in real time.

All these examples involve interaction between an agent & its environment; where the agent seeks to achieve a goal despite uncertainty about its environment (next move at 'go', the robot's next location, where to invest (which assets, how much), levels of the reservoirs). The agent's actions affect the options available at later times ( $\rightarrow$  feedback is delayed, not instantaneous). Correct choices require taking into account delayed consequences of actions  $\rightarrow$  planning required; think ahead!

## Elements of RL

2

- A policy = how the agent picks its actions.  
= a mapping from perceived states of the environment to actions to be taken.
- A reward signal = how well an agent is doing at time  $t$ .  
The agent's goal is to maximize the cumulative reward.  
 $\equiv$  immediate reward  
 $\rightarrow$  primary basis for altering the policy.
- A value function = captures what is good in the long run (as opposed to the reward signal): how good is each state and/or actions?  
The value of a state is the total amount of reward an agent can expect to accumulate over the future, starting from that state.
- A model of the environment = the agent's representation of the environment.  
 $\rightarrow$  Optional. Methods in RL that use a model are known as model-based approaches (first learn the environment model, and then use that to learn a policy), as opposed to model-free approaches (learning an action policy directly).

### I. FINITE MARKOV DECISION PROCESSES.

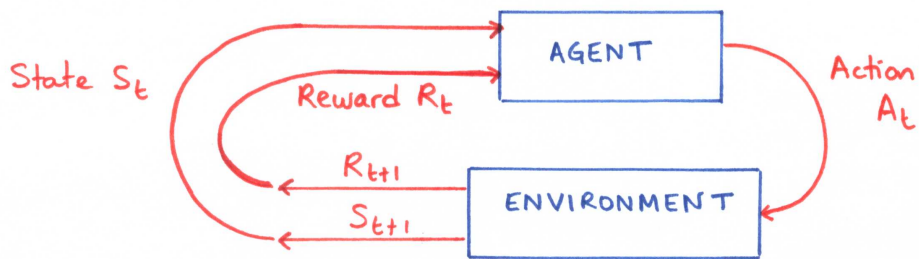
Markov Decision Processes (MDP) are ideal mathematical models for the RL problem.

$\rightarrow$  The learner is called agent

$\rightarrow$  Everything outside the agent is called the environment.

The agent & the environment continuously interact: at each step, the agent executes an action  $A_t$ , given some information about

the environment's state  $S_t$ . As a consequence of its action, (3) the agent receives a reward  $R_{t+1}$ , and finds itself in a new state  $S_{t+1}$ .



⇒ gives rise to a trajectory:  $S_0, A_0, R_1, S_1, A_1, R_2, S_2, A_2, \dots$

We assume a finite MDP = states, actions, and rewards can only take finitely many values:

$$\begin{aligned}
 S_t &\in \mathcal{S} \\
 A_t &\in \mathcal{A} \\
 R_t &\in \mathcal{R}
 \end{aligned}$$

sets  $\mathcal{S}, \mathcal{A}, \mathcal{R}$  contain finitely many elements.

→ Dynamics of the MDP: Random Variables  $R_t, S_t$  have discrete probability distributions that depend only on the preceding state & action (the Markov property):

$$p(s', r | s, a) := \mathbb{P}(S_{t+1} = s', R_{t+1} = r | S_t = s, A_t = a)$$

$$\forall s, s' \in \mathcal{S}, \forall r \in \mathcal{R}, \forall a \in \mathcal{A}.$$

sum to one.

Remark: the state representation is highly important to determine what happens next. If the state representation is not complete enough, then the Markov property is no longer justified, and the techniques developed may fail.

Ex: trajectory of a helicopter. The state variable should incorporate position, angular position, angular velocity, wind

direction, and so on. A state representation including position (4) only is not sufficient to control the helicopter and choose actions that will lead it to some target: we would need to go back in time and consider its full trajectory (→ Markov property violated).

From the function  $p$ , we can compute important quantities such as

\* state-transition probabilities.

$$\begin{aligned}
 p(s' | s, a) &:= \mathbb{P}(S_{t+1} = s' | S_t = s, A_t = a) \\
 &= \sum_{r \in \mathcal{R}} p(s', r | s, a).
 \end{aligned}$$

\* expected rewards for (state, action) pairs

$$\begin{aligned}
 r(s, a) &:= \mathbb{E}(R_{t+1} | S_t = s, A_t = a) \\
 &= \sum_{r \in \mathcal{R}} r \sum_{s' \in \mathcal{S}} p(s', r | s, a).
 \end{aligned}$$

→ The sequence of rewards received after time  $t$  is  $R_{t+1}, R_{t+2}, \dots$ . We seek to maximize the (expected) return  $\mathbb{E} G_t$ , where

$$G_t = \underbrace{R_{t+1} + R_{t+2} + \dots + R_T}_{\text{cumulative reward}}$$

finite horizon  
finite time step.

In many cases, it is convenient to consider infinite horizon:  $T = +\infty$ . [Ex: robot with long life span]. The return formulation is problematic, since  $G_t$  could easily be infinite.

⇒ We need the concept of discounting.

The discounted return  $\tilde{w}$ :  $G_t := R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots$

$$= \sum_{k \geq 0} \gamma^k R_{t+k+1}, \quad \gamma \in [0, 1]$$

$\gamma =$  discount rate.

Why introducing  $\gamma$ ?

(5)

(i) Mathematically convenient: if  $\gamma < 1$ , the infinite sum converges provided that the reward sequence is bounded.

(ii) Trade-off between immediate reward (extreme case:  $\gamma = 0$ : "myopic" strategy as the agent is concerned only by its immediate rewards) and "far-sighted" objectives (as  $\gamma$  gets close to 1, the return  $G_t$  takes future rewards into account more strongly).

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

$$\Rightarrow G_t = R_{t+1} + \gamma G_{t+1} \quad (*)$$

↑ This relation will have strong implications later.

To summary, a MDP is a tuple  $(\mathcal{S}, \mathcal{A}, \mathcal{T}, r, \gamma)$

transitions  $p(s'|s, a)$       discount rate  $\gamma$   
 expected rewards  $r(s, a)$

• A policy is a mapping from  $\mathcal{S}$  to  $\mathcal{A}$ :  $\pi: \mathcal{S} \rightarrow \mathcal{A}$

• deterministic policy  $a = \pi(s)$

• stochastic policy  $\pi(a|s) = \mathbb{P}(A=a | S=s)$   
 (stationary policy: does not depend on  $t$ )

↑  
 The agent's behaviour function.

• The value function  $v_\pi$  of a state  $s$  given a policy  $\pi$  quantifies the expected return when starting in  $s$  & following  $\pi$ .

$$v_\pi(s) := \mathbb{E}_\pi(G_t | S_t = s)$$

↑ ↑  
 Emphasizes the dependence of the value function on a particular  $\pi$

• Another important quantity is the action-value function  $q_\pi$ , (6) which characterizes the expected return when starting in  $s$ , taking action  $a$ , and following policy  $\pi$ :

$$q_\pi(s, a) := \mathbb{E}_\pi(G_t | S_t = s, A_t = a)$$

• A fundamental property of the value function & action-value function is that they can be decomposed into two parts:  
 - immediate reward  
 + discounted value of future reward:

From the decomposition (\*), we have

$$\begin{aligned} v_\pi(s) &= \mathbb{E}_\pi(R_{t+1} + \gamma G_{t+1} | S_t = s) \\ &= \mathbb{E}_\pi(R_{t+1} | S_t = s) + \gamma \mathbb{E}_\pi(G_{t+1} | S_t = s) \\ &= \sum_{a \in \mathcal{A}} \pi(a|s) r(s, a) + \gamma \mathbb{E}_\pi(G_{t+1} | S_t = s) \\ &= \mathbb{E}_\pi\{r(s, A_t) | S_t = s\} + \gamma \mathbb{E}_\pi\{v_\pi(S_{t+1}) | S_t = s\} \end{aligned}$$

Since

$$\mathbb{E}_\pi(R_{t+1} | S_t = s) = \mathbb{E}\{E(R_{t+1} | S_t, A_t) | S_t = s\} = \mathbb{E}\{r(S_t, A_t) | S_t = s\}$$

$$\begin{aligned} &\mathbb{E}\{E(v_\pi(S_{t+1}) | S_t = s, A_t = a) | S_t = s\} \\ &= \text{function of } s, a \\ &= \sum_{s' \in \mathcal{S}} v_\pi(s') p(s'|s, a) \\ &= \text{function of } s = \sum_{a \in \mathcal{A}} \pi(a|s) \sum_{s' \in \mathcal{S}} v_\pi(s') p(s'|s, a) \end{aligned}$$

We obtain the decomposition:

7

$$V_{\pi}(s) = \sum_{a \in A} \pi(a|s) \left\{ r(s, a) + \gamma \sum_{s' \in S} p(s'|s, a) V_{\pi}(s') \right\}$$

aka the BELLMAN EQUATION for  $V_{\pi}$ .

With a deterministic policy, Bellman equation reduces to

$$V_{\pi}(s) = r(s, \pi(s)) + \gamma \sum_{s' \in S} p(s'|s, \pi(s)) V_{\pi}(s')$$

In matrix notation:

$$V_{\pi} = R_{\pi} + \gamma P_{\pi} V_{\pi}$$

$$V_{\pi} := \begin{pmatrix} | \\ V_{\pi}(s) \\ | \end{pmatrix}$$

$$R_{\pi} := \begin{pmatrix} | \\ r(s, \pi(s)) \\ | \end{pmatrix}$$

$$P_{\pi} := \begin{pmatrix} | & & | \\ \hline p(s'|s, \pi(s)) & & \hline \end{pmatrix} \begin{matrix} \leftarrow \\ \leftarrow \\ \leftarrow \end{matrix} = |S| \times |S| \text{ transition matrix} \\ \begin{matrix} \uparrow \\ \uparrow \\ \uparrow \end{matrix} \begin{matrix} s\text{-th row} \\ \\ \\ \end{matrix} = (P_{s, s'}) \\ \begin{matrix} \uparrow \\ \uparrow \\ \uparrow \end{matrix} \begin{matrix} s'\text{-th column} \\ \\ \\ \end{matrix}$$

The (finite) system of linear equations  $V_{\pi} = R_{\pi} + \gamma P_{\pi} V_{\pi}$  admits a unique solution  $V_{\pi}^* = (I - \gamma P_{\pi})^{-1} R_{\pi}$ . To prove this, we only have to verify that  $I - \gamma P_{\pi}$  is invertible.

$$\begin{aligned} \|P_{\pi}\|_{\infty} &= \max_s \sum_{s'} |P_{s, s'}| \\ &= \max_s \sum_{s'} \mathbb{P}(S_{t+1} = s' | S_t = s, A_t = \pi(s)) \\ &= 1 \end{aligned}$$

$\Rightarrow \| \gamma P_{\pi} \|_{\infty} = \gamma < 1 \Rightarrow$  The eigenvalues of  $\gamma P_{\pi}$  are all less than 1, and  $(I - \gamma P_{\pi})$  is thus invertible.

• Similarly,

8

$$\begin{aligned} q_{\pi}(s, a) &= \mathbb{E}_{\pi}(G_t | S_t = s, A_t = a) \\ &= \mathbb{E}_{\pi}(R_{t+1} + \gamma G_{t+1} | S_t = s, A_t = a) \\ &= \mathbb{E}_{\pi}(R_{t+1} + \gamma q_{\pi}(S_{t+1}, A_{t+1}) | S_t = s, A_t = a) \\ &= r(s, a) + \gamma \underbrace{\mathbb{E}_{\pi}\{q_{\pi}(S_{t+1}, A_{t+1}) | S_t = s, A_t = a\}} \\ &\quad \parallel \\ &\quad \sum_{s', a'} q(s', a') \mathbb{P}(S_{t+1} = s', A_{t+1} = a' | S_t = s, A_t = a) \\ &\quad \parallel \\ &\quad \sum_{s', a'} q(s', a') \mathbb{P}(A_{t+1} = a' | S_t = s, A_t = a, S_{t+1} = s') \\ &\quad \quad \quad \parallel \\ &\quad \quad \quad \equiv \pi(a' | s') \\ &\quad \quad \quad \times \mathbb{P}(S_{t+1} = s' | S_t = s, A_t = a) \\ &\quad \quad \quad \parallel \\ &\quad \quad \quad \equiv p(s' | s, a) \\ &\quad \quad \quad \parallel \\ &\quad \quad \quad \sum_{s', a'} q(s', a') \pi(a' | s') p(s' | s, a) \\ &\quad \quad \quad \parallel \\ &\quad \quad \quad \sum_{s'} p(s' | s, a) \underbrace{\sum_{a'} \pi(a' | s') q(s', a')} \\ &\quad \quad \quad \parallel \\ &\quad \quad \quad \equiv V_{\pi}(s') \end{aligned}$$

Since

$$\begin{aligned} V_{\pi}(S_t) &= \mathbb{E}_{\pi}\{G_t | S_t\} \\ &= \mathbb{E}_{\pi}\{\mathbb{E}_{\pi}(G_t | S_t, A_t) | S_t\} \\ &= \mathbb{E}_{\pi}\{q_{\pi}(S_t, A_t) | S_t\} \\ \Rightarrow V_{\pi}(s') &= \sum_{a' \in A} \pi(a' | s') q_{\pi}(s', a'). \end{aligned}$$

We finally obtain

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

BELLMAN EQUATION for  $q_{\pi}$ .

• Optimal policies & Optimal Value Functions.

→ A policy  $\pi$  is said to perform better than  $\pi'$  if  $v_\pi(s) \geq v_{\pi'}(s) \quad \forall s \in \mathcal{S}$ ; we write  $\pi \geq \pi'$ .

There is always a policy that is better than or equal to all other policies, but it may not be unique. If there is more than one, we denote by  $\pi_*$  all optimal policies. They all share the same value function  $v_* : v_*(s) := \max_{\pi} v_\pi(s)$ .

We see from Bellman equation for  $q_\pi$  that all optimal policies  $\pi_*$  also share the same optimal action-value function

$$\begin{aligned}
 q_*(s, a) &= r(s, a) + \gamma \sum_{s'} p(s'|s, a) v_*(s') \\
 &= r(s, a) + \gamma \sum_{s'} p(s'|s, a) \max_{\pi} v_\pi(s) \\
 &= \max_{\pi} \left\{ r(s, a) + \gamma \sum_{s'} p(s'|s, a) v_\pi(s) \right\} \\
 &= \max_{\pi} q_\pi(s, a)
 \end{aligned}$$

Next, note that  $v_*(s) = \sum_{a \in \mathcal{A}} \pi_*(a|s) \left\{ r(s, a) + \gamma \sum_{s' \in \mathcal{S}} p(s'|s, a) v_*(s') \right\}$

Bellman equation for  $v_*$

$$\begin{aligned}
 &= \sum_{a \in \mathcal{A}} \pi_*(a|s) q_*(s, a) \\
 &\leq \sum_{a \in \mathcal{A}} \pi_*(a|s) \left( \max_a q_*(s, a) \right) \\
 &= \max_a q_*(s, a).
 \end{aligned}$$

In addition, if there exists a state  $s$  for which this inequality is strict:  $v_*(s) < \max_a q_*(s, a)$ , then taking a deterministic policy with the maximizing action would define a better policy  $\Rightarrow$  we must have equality:  $v_*(s) = \max_a q_*(s, a) \quad \forall s$ :

$$v_*(s) = \max_a \left\{ r(s, a) + \gamma \sum_{s'} p(s'|s, a) v_*(s') \right\}$$

← BELLMAN EQUALITY EQUATION for  $v_*$ .

In addition, we observe that

$$\forall s \in \mathcal{S}, \quad \pi_*(s) = \operatorname{argmax}_a q_*(s, a)$$

If we know  $q_*$ , we immediately know  $\pi_*$  (and  $v_*$ )

x Remark: Bellman optimality equation for  $q_*$  is

$$q_*(s, a) = r(s, a) + \gamma \sum_{s'} p(s'|s, a) \max_{a'} q_*(s', a')$$

→ Bellman optimality equations are non-linear due to the presence of the max  $\rightarrow$  in general, no closed form solution. However, many iterative solutions exist, and are presented in the next section.

x Remark: The knowledge of  $q_*$  make the choice of an optimal policy very easy. However, knowing  $v_*$  also allow us to deduce an optimal policy: for each  $s$ , there will be one or more actions  $a$  that achieve the maximum on the RHS of Bellman optimality equations for  $v_*$   $\Rightarrow$  any policy that assigns a non-zero probability to these actions only is an optimal policy.

Since for these values of  $a$ , Bellman equation is indeed satisfied:

$$v_*(s) = \underbrace{\sum_{a \in \mathcal{A}} \pi_*(a|s)}_{\text{sum to 1}} \underbrace{\left\{ r(s, a) + \gamma \sum_{s' \in \mathcal{S}} p(s'|s, a) v_*(s') \right\}}_{\substack{\parallel \\ \max_a \{ \dots \} = v_*(s)}}$$

since  $\pi_*$  puts non-zero mass only on the values of  $a$  reaching the max.

## II - DYNAMIC PROGRAMMING.

(11)

Dynamic Programming (DP) refers to a family of algorithms that can be used to compute optimal policies given the MDP parameters. (term DP is due to Bellman '57)

↳ No learning:  $p(s', r | s, a)$  is given.

⇒ A PLANNING problem.

We discuss in this section two planning algorithms: policy iteration & value iteration.

### II.1. Policy Evaluation.

We first address a simple problem: given a policy  $\pi$  (deterministic or stochastic), how to compute the state-value function  $v_\pi(s)$ ,  $\forall s \in \mathcal{S}$ ?

→ Recall Bellman equation from page 7:

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

(\*)

$$V_\pi = R_\pi + \gamma P_\pi V_\pi \quad (\text{matrix form})$$

$$\bullet V_\pi = \begin{pmatrix} | & & | \\ v_\pi(s) & & \\ | & & | \end{pmatrix} \leftarrow \begin{matrix} \text{s-th} \\ \text{row} \end{matrix} \bullet P_\pi = \begin{pmatrix} \text{---} & \sum_{a \in \mathcal{A}} \pi(a|s) p(s'|s, a) & \text{---} \end{pmatrix} \leftarrow \begin{matrix} \text{s-th} \\ \text{column} \end{matrix}$$

$$\bullet R_\pi = \begin{pmatrix} \text{---} & \sum_{a \in \mathcal{A}} \pi(a|s) r(s, a) & \text{---} \end{pmatrix} \leftarrow \begin{matrix} \text{s-th} \\ \text{row} \end{matrix}$$

↑  
s'-th column

(\*) = system of  $|\mathcal{S}|$  linear equations with  $|\mathcal{S}|$  unknowns.

(12)

It can be solved

↳ using matrix inversion:  $V_\pi = (I - \gamma P_\pi)^{-1} R_\pi$

↳ iteratively:  $V_\pi^{(k+1)} = R_\pi + \gamma P_\pi V_\pi^{(k)}$



$$\begin{aligned} \text{Note that } \|P_\pi\|_\infty &= \max_s \sum_{s'} |p(s, s')| \\ &= \max_s \sum_{s'} \sum_a \pi(a|s) p(s'|s, a) \\ &= \max_s \sum_a \pi(a|s) \underbrace{\sum_{s'} p(s'|s, a)}_{=1} = 1 \end{aligned}$$

⇒  $\|\gamma P_\pi\|_\infty = \gamma < 1 \Rightarrow$  Eigenvalues of  $\gamma P_\pi$  are all  $< 1$ .

This ensures that (i) the matrix  $(I - \gamma P_\pi)$  is invertible  
(ii) sequence  $V_\pi^{(k)}$  converges to the solution.

• Direct inversion of  $(I - \gamma P_\pi)$  may be tedious & computationally expensive, and we use the iteration method, more suitable for our purposes:

#### ITERATIVE POLICY EVALUATION

$\forall s \in \mathcal{S}$ , compute (until  $\max_s |v_\pi^{(k+1)}(s) - v_\pi^{(k)}(s)| < \epsilon$ )

$$v_\pi^{(k+1)}(s) = \sum_{a \in \mathcal{A}} \pi(a|s) \left\{ r(s, a) + \gamma \sum_{s' \in \mathcal{S}} p(s'|s, a) v_\pi^{(k)}(s') \right\}$$

↑  
All states are updated using the old value  $v_\pi^{(k)}$ . aka SYNCHRONOUS UPDATES. Alternatively, after one state  $s$  is updated, the new value function may be plugged back into the right-hand side, overwriting the old one → asynchronous updates.

Ex: Consider the following 4x4 grid:

	1	2	3
4	5	6	7
8	9	10	11
12	13	14	

terminal states:  
if you arrive here; you stay here

- States: Where on the grid you are.
- Actions: Move North, South, East or West only  $\left\langle \begin{array}{c} \uparrow \\ \leftarrow \\ \downarrow \\ \rightarrow \end{array} \right\rangle$   
(from 6, you can only get to 2, 5, 7, 10)  
(from 7, you can get to 3, 6, 11, or remain at 7, shall you decide to move East).  
(from 3, you can go to 2 or 7, or remain at 3 shall you decide to go North or East).

- Rewards:  $R_t = -1$  on all transitions.
- Policy: Uniform policy: moves N, S, E or W occur with probability  $1/4$ :  $\pi(a | s) = 1/4$  for  $a \in \{N, S, E, W\}$ .

Transition probabilities are of the form:

$$p(s', r | s, a) = 1$$

$\begin{matrix} \uparrow & \uparrow & \uparrow & \uparrow \\ s' & r & s & a \end{matrix}$ 
 (notation from page 3)

$$p(6, r | 13, \text{south}) = 0 \quad \forall r \quad \dots / \dots$$

Goal: Evaluate  $v_\pi(s)$  for all  $s$  (i.e. for all positions on the grid)

Use Iterative Policy Evaluation Algorithm.

Put  $R_\pi := \begin{pmatrix} 0 \\ -1 \\ \vdots \\ -1 \end{pmatrix}$   $\left\{ \begin{array}{l} \leftarrow \text{terminal state } \blacksquare = \text{no reward} \\ \leftarrow 14 \text{ entries (fixed rewards } = -1) \\ \leftarrow \text{terminal state } \blacksquare = \text{no reward} \end{array} \right.$

		to state s'															
from state s		1	2	3	4	5	6	7	8	9	10	11	12	13	14		
$P_\pi =$	1	1/4	1/4	1/4			1/4										
	2		1/4	1/4	1/4				1/4								
	3			1/4	1/2					1/4							
	4	1/4			1/4	1/4					1/4						
	5		1/4		1/4	1/4		1/4			1/4						
	6			1/4		1/4		1/4		1/4			1/4				
	7				1/4		1/4	1/4						1/4			
	8					1/4			1/4	1/4					1/4		
	9						1/4			1/4	1/4					1/4	
	10							1/4			1/4		1/4				1/4
	11								1/4			1/4	1/4				1/4
	12									1/4				1/2	1/4		
	13										1/4			1/4	1/4	1/4	
	14											1/4			1/4	1/4	1/4

= 0 everywhere else

Each entry is equal to  $\sum_{a \in A} \underbrace{\pi(a|s)}_{1/4} \underbrace{p(s'|s, a)}_{\text{all 0/1 proba}}$

Take  $V_\pi^{(0)} = \begin{pmatrix} 0 \\ 1 \\ 0 \end{pmatrix}$   $16 \times 1$  & put  $V_\pi^{(k+1)} = R_\pi + \gamma P_\pi V_\pi^{(k)}$   $k=0, 1, \dots$

0	-1	-1	-1
-1	-1	-1	-1
-1	-1	-1	-1
-1	-1	-1	0

$k=1$

0	-1.75	-2	-2
-1.75	-2	-2	-2
-2	-2	-2	-1.75
-2	-2	-1.75	0

$k=2$

0	-2.44	-2.94	-3
-2.44	-2.87	-3	-2.94
-2.94	-3	-2.87	-2.44
-3	-2.94	-2.44	0

$k=3$

Value of the uniform policy after k iterations

0	-3.66	-4.69	-4.91	0	-6.14	-8.35	-8.97	0	-14	-20	-22
-3.66	-4.48	-4.78	-4.69	-6.14	-7.74	-8.43	-8.35	-14	-18	-20	-20
-4.69	-4.78	-4.48	-3.66	-8.35	-8.43	-7.74	-6.14	-20	-20	-18	-14
-4.91	-4.69	-3.66	0	-8.97	-8.35	-6.14	0	-22	-20	-14	0

$k=5$

$k=10$

$k=\infty$   
(solution  $V_\pi = (I - \gamma P_\pi)^{-1} R_\pi$ )

Interpretation: each number in each cell represents the average number of steps starting from that cell, before reaching a terminal absorbing state.

## II.2. Policy Improvement.

Suppose we have determined the value  $v_\pi$  of a policy  $\pi$ .  
How to change its value at some state  $s$  to find a better policy

• Tool: the action-value function. Recall:

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

Interpretation: value of taking action  $a$  at state  $s$ , and the following policy  $\pi$ .  
→ If you find a value  $a$  such that  $q_\pi(s, a)$  is larger than  $v_\pi(s)$ , i.e. if it is better to select action  $a$  once from state  $s$  & then

follow policy  $\pi$ , then you would expect that you are better off following action  $a$  from  $s$  all the time. And indeed, we have the following result:

Theorem. Consider a deterministic policy  $\pi, \pi'$

$$\text{If } \forall s \in \mathcal{S} \quad q_\pi(s, \pi'(s)) \geq v_\pi(s) \quad (1)$$

Then

$$\forall s \in \mathcal{S} \quad v_{\pi'}(s) \geq v_\pi(s) \quad (2)$$

In addition, if there is at least one strict inequality at (1), the inequality at (2) is also strict.

proof:

$$\begin{aligned} v_\pi(s) &\leq q_\pi(s, \pi'(s)) \\ &= r(s, \pi'(s)) + \gamma \sum_{s'} p(s' | s, \pi'(s)) \underbrace{v_\pi(s')}_{\leq q_\pi(s', \pi'(s'))} \\ &\leq r(s, \pi'(s)) + \gamma \sum_{s'} p(s' | s, \pi'(s)) \underbrace{q_\pi(s', \pi'(s'))}_{=} \\ &\quad r(s', \pi'(s')) + \gamma \sum_{s''} p(s'' | s', \pi'(s')) v_\pi(s'') \\ &= r(s, \pi'(s)) + \gamma \sum_{s'} p(s' | s, \pi'(s)) r(s', \pi'(s')) \\ &\quad + \gamma^2 \sum_{s', s''} p(s' | s, \pi'(s)) p(s'' | s', \pi'(s')) \underbrace{v_\pi(s'')}_{\leq q_\pi(s'', \pi'(s''))} \\ &\dots / \dots \text{ keep iterating} \\ &\leq v_{\pi'}(s) (= \mathbb{E}(R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots | S_t = s)) \end{aligned}$$

15

16



→ Instead of changing a single state  $s$  only, it is natural to consider changes at all states, and consider the new GREEDY policy  $\pi'$ : (17)

POLICY IMPROVEMENT

$$\forall s \in \mathcal{S}$$

$$\begin{aligned} \pi'(s) &= \operatorname{argmax}_a q_{\pi}(s, a) \\ &= \operatorname{argmax}_a \left\{ r(s, a) + \gamma \sum_{s'} p(s'|s, a) v_{\pi}(s') \right\} \end{aligned}$$

→ Suppose now that the new policy  $\pi'$  is as good as  $\pi$ , but not strictly better:  $v_{\pi'}(s) = v_{\pi}(s) \forall s$ . Then we have

$$v_{\pi'}(s) = \max_a \left\{ r(s, a) + \gamma \sum_{s'} p(s'|s, a) v_{\pi'}(s') \right\}$$

↑ ( $v_{\pi''}(s')$ )

This is exactly Bellman inequality equation.

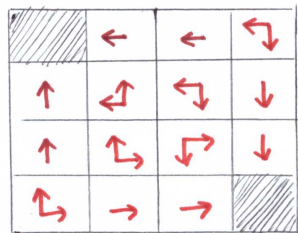
⇒  $v_{\pi'} = v_{\pi} = v^*$  & both  $\pi/\pi'$  are optimal policies.

⇒ Policy improvement is guaranteed to give us a strictly better policy, unless the optimal policy is reached.

Ex: Back to the grid example on page 13. The uniform policy has value

0	-14	-20	-22
-14	-18	-20	-20
-20	-20	-18	-14
-22	-20	-14	0

greedy  
policy  
improvement



policy  $\pi'$ .

In this example, the improved policy  $\pi'$  turns out to be optimal. This is not true in general. (18)

II.3. Policy Iteration.

The policy iteration algorithm alternates between policy evaluation & greedy policy improvement, generating a sequence

$$\pi_0 \rightarrow v_{\pi_0} \rightarrow \pi_1 \rightarrow v_{\pi_1} \rightarrow \pi_2 \rightarrow v_{\pi_2} \rightarrow \dots$$

Each iteration must be a strict improvement of the policy, unless we are at the optimal. Since a finite MDP has finitely many policies, the sequence must converge to the optimal policy in a finite number of steps.

• Goal: compute the optimal policy  $\pi^*$ .

• Init: arbitrary policy  $\pi_0$

$$\begin{aligned} \pi &\leftarrow \pi_0 \\ \pi' &\leftarrow \text{NIL} \end{aligned}$$

• While ( $\pi \neq \pi'$ ) do

    ↳ Policy Evaluation: determine  $v_{\pi}$  (page 12);  $V \leftarrow v_{\pi}$

$$\pi' \leftarrow \pi$$

    ↳ Greedy Policy Improvement:

$$\pi = \operatorname{argmax}_{\pi} \{ R_{\pi} + \gamma P_{\pi} V \}$$

• Return  $\pi$

POLICY ITERATION

## II. 4. Value Iteration.

(19)

- Given a policy  $\pi$  with value function  $v_{\pi}(s)$

$$r(s, \pi(s)) + \gamma \sum_{s'} p(s'|s, \pi(s)) v_{\pi}(s')$$

we obtain a better policy  $\pi'$  by selecting for each  $s$ , the value a maximizing  $r(s, a) + \gamma \sum_{s'} p(s'|s, a) v_{\pi}(s')$ .

↳ Back to the grid example, the value function of the random policy  $\frac{1}{4} \begin{matrix} \leftarrow \\ \uparrow \\ \pi \\ \downarrow \\ \rightarrow \end{matrix} \frac{1}{4}$  can be used to greedily select a

better policy  $\pi'$  (which is optimal in this particular example). However, if the value function  $v_{\pi}$  is calculated iteratively, we see that a greedy policy based on the third iterate ( $k=3$ ) yields already the same optimal policy  $\pi'$  (see page 14). This observation shows that we may not need to wait until convergence of the iterative policy evaluation algorithm to greedily select the next policy; which raises the question: after how many iterations stop?

- The value iteration algorithm stops after exactly one iteration during the evaluation step:

→ Given  $v^{(0)}(s)$  [initialization: arbitrary value]

→ Calculate (policy improvement algo with early stopping)

$$v^{(1)}(s) = \max_{\pi} \left\{ r(s, \pi(s)) + \gamma \sum_{s'} p(s'|s, \pi(s)) v^{(0)}(s') \right\}$$

deterministic policy

1. Evaluation step (stop after 1<sup>st</sup> iteration)

2. Improvement step: greedy search.

& repeat:

(20)

$$v^{(k+1)}(s) = \max_{\pi} \left\{ r(s, \pi(s)) + \gamma \sum_{s'} p(s'|s, \pi(s)) v^{(k)}(s') \right\}$$

↑ Another way to understand this expression is by reference to Bellman optimality equation for  $v^*$  (see page 9): the Bellman equation is simply made recursive.

↳ In matrix notation: 
$$V^{(k+1)} = \max_{\pi} \left\{ R_{\pi} + \gamma P_{\pi} V^{(k)} \right\}$$
  

$$=: \Phi(V^{(k)})$$

**Theorem:** For any initial value  $V^{(0)}$ , the sequence of iterates  $V^{(k+1)} = \Phi(V^{(k)})$  converges to the optimal value function  $V^*$ .

proof: We show that  $\Phi$  is  $\gamma$ -Lipschitz for  $\|\cdot\|_{\infty}$ . Since  $\gamma < 1$ ,  $\Phi$  is then a contraction & admits a fixed point; and any sequence of iterates converges to the fixed point.

Let  $s \in \mathcal{S}$ ,  $V \in \mathbb{R}^{|\mathcal{S}|}$ ,  $U \in \mathbb{R}^{|\mathcal{S}|}$ ,

& denote by  $a^*(s) = \operatorname{argmax}_a \left\{ r(s, a) + \gamma \sum_{s'} p(s'|s, a) v(s') \right\}$

Then

$$\begin{aligned} \Phi(V)(s) - \Phi(U)(s) &\leq \Phi(V)(s) - \left\{ r(s, a^*(s)) + \gamma \sum_{s'} p(s'|s, a^*(s)) U(s') \right\} \\ &= \gamma \sum_{s'} p(s'|s, a^*(s)) \{ V(s') - U(s') \} \\ &\leq \gamma \sum_{s'} p(s'|s, a^*(s)) \|V - U\|_{\infty} \\ &= \gamma \|V - U\|_{\infty}. \end{aligned}$$

s-th entry of the vector  $\Phi(V)$

& similarly  $\Phi(u)(s) - \Phi(v)(s) \leq \gamma \|v - u\|_\infty$ , (21)  
 so that  $|\Phi(v)(s) - \Phi(u)(s)| \leq \gamma \|v - u\|_\infty$ .

$\Phi$  is  $\gamma$ -Lipschitz.

Thus

$$\begin{aligned} \|v^* - v^{(k+1)}\|_\infty &= \|\Phi(v^*) - \Phi(v^{(k)})\|_\infty \\ &\leq \gamma \|v^* - v^{(k)}\|_\infty \\ &\leq \gamma^{k+1} \|v^* - v^{(0)}\|_\infty \\ &\rightarrow 0 \text{ as } k \rightarrow +\infty. \end{aligned}$$

- Remarks
- (i) Intermediate value functions may not correspond to any policy: it can happen that  $\nexists \pi$  such that  $v^{(k)} = \mathcal{J}_\pi v^{(k)}$ .
  - (ii) The algorithm  $v^{(k+1)} = \Phi(v^{(k)})$  uses synchronous backups. Instead, we may use asynchronous backups.
  - (iii) Policy & Value Iteration algorithms are based on state-value functions. The complexity is  $O(m|S|^2)$  per iteration, for  $m$  actions and  $|S|$  states. We could apply the same ideas to the action-value functions. Complexity is  $O(m^2|S|^2)$ .  
 For each state, we need to consider the successor state ( $\rightarrow |S|^2$  operations), and there are  $m$  actions.

- Goal: compute the optimal policy  $\pi^*$  / value function  $v^*$
- Init:  $v^{(0)}$
- Repeat  $k=0, 1, \dots$   

$$v^{(k+1)} = \max_{\pi} \{ R_{\pi} + \gamma P_{\pi} v^{(k)} \}$$
- Until Convergence.

VALUE ITERATION

### III - MODEL-FREE PREDICTION

MDP environment but we do not know the model parameters (i.e. transition & reward probabilities).

$\hookrightarrow$  We want to find the optimal way to behave.

In this section, we give methods to estimate the value function of a given policy, when the environment is unknown.

In section IV, we optimize the value function of an unknown MDP (control).

We discuss two classes of techniques:

- MC learning
- TD learning.

These methods all learn from experience: given a policy  $\pi$ , sample trajectories according to  $\pi$ , observe rewards & successive states sampled by the environment, and update the value function estimate based on sample returns.

$\uparrow$  learning from actual experience

#### III.1. Monte Carlo Methods.

- Assumption: Experience is divided into episodes, and all episodes eventually terminate, no matter what actions are performed.

Episode #1 :  $S_0^{(1)}, A_0^{(1)}, R_1^{(1)}, S_1^{(1)}, A_1^{(1)}, R_2^{(1)} - \text{end}$   
 Episode #2 :  $S_0^{(2)}, A_0^{(2)}, R_1^{(2)} - \text{end}$   
 Episode #3 :  $S_0^{(3)}, A_0^{(3)}, R_1^{(3)}, S_1^{(3)}, A_1^{(3)}, R_2^{(3)} - \text{end}$   
 ...

At the end of an episode, the value estimates are updated (episode-to-episode updates; not "on-line").  
 Actions are selected by the agent; states & rewards are received by the environment (no need to know its dynamics)

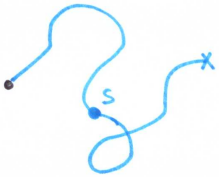
Remark: Learning is from complete episodes only.

(23)

- Suppose we wish to estimate the value  $v_{\pi}(s)$  of policy  $\pi$  at state  $s$ .

Let  $t =$  the first time in an episode that state  $s$  is visited ( $s$  can be visited multiple times).

Ex:  $s' \rightarrow a' \rightarrow r' \rightarrow \textcircled{s} \rightarrow a'' \rightarrow r'' \rightarrow s' \rightarrow a' \rightarrow r'$   
 $\rightarrow \textcircled{s} \rightarrow a'' \rightarrow r''' - \text{end}$   
first visit to  $s$   
second visit to  $s$ .



Return at time  $t$  is  $G_t = R_{t+1} + \gamma R_{t+2} + \dots + \gamma^{T-t} R_T$ .

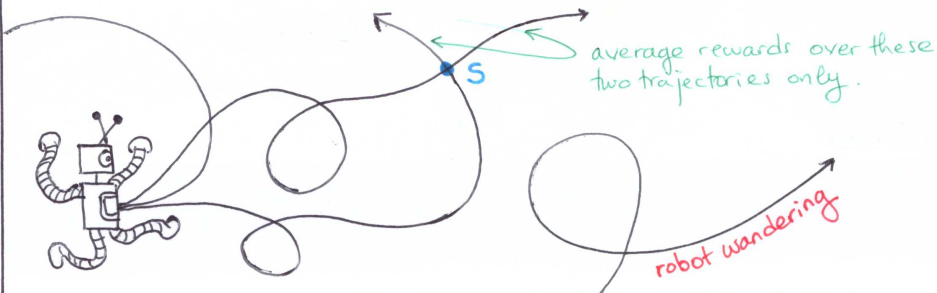
& the value function at state  $s$  is

$$v_{\pi}(s) = \mathbb{E}(G_t | S_t = s)$$

Monte-Carlo (MC) techniques use the empirical mean returns instead of the population mean. (SLLN ensure convergence)

• 2 paradigms:

↳ "First-visit" MC policy prediction considers the reward  $G_t$  associated with the first-time step  $t$  that state  $s$  is visited. The value  $v_{\pi}(s)$  is estimated by sample mean return.



Algorithm: First-visit MC prediction

(24)

x Goal: Estimate  $v_{\pi}$  of a policy  $\pi$  (unknown environment)

x Init:  $v(s) \in \mathbb{R}$  arbitrary  $\forall s \in \mathcal{S}$ .

• Returns  $\leftarrow$  empty list  $\forall s \in \mathcal{S}$ .

x Repeat (for each episode)

(1) Generate episode following policy  $\pi$ :

$S_0, A_0, R_1, S_1, A_1, R_2, \dots, S_{T-1}, A_{T-1}, R_T$ .

(2)  $G \leftarrow 0$

(3) For  $t = T-1, T-2, \dots, 0$ :

$G \leftarrow \gamma G + R_{t+1}$

Unless  $S_t$  appears in  $S_0, S_1, \dots, S_{t-1}$ :

Add  $G$  to Returns ( $S_t$ )

Update  $v(S_t) \leftarrow$  Average (Returns ( $S_t$ )).

↳ "Every-visit" MC policy prediction averages the returns following every visit to  $s$ . The "every-visit" algorithm is the same as the algorithm presented above, but without the line "Unless  $S_t$  appears in  $S_0, \dots, S_{t-1}$ ".

• Result: Both first-visit & every-visit MC converge to  $v_{\pi}(s)$  as the number of visits to  $s$  goes to infinity.

First-visit  $\rightarrow$  A direct consequence of the SLLN (a.s. convergence)

Every-visit  $\rightarrow$  see Singh & Sutton (1996)

• Caveat: We have to wait until the end of an episode to compute the mean. We present next TD learning, which learns from incomplete episodes, by bootstrapping.

### III.2. TD learning.

(25)

- Temporal-Difference (TD) learning is a combination of MC ideas & Dynamic Programming (DP) ideas.
  - MC: learn from experience without knowing the environment's dynamics.
  - DP: update estimates from already learned estimates, without waiting for the final outcome.

In MC learning, the value of a state  $s$  is estimated by averaging the returns  $G_t$ :

"Update  $v(S_t) \leftarrow \text{Average}(\text{Returns}(S_t))$ ."

Instead of recomputing the sample mean at each iteration, the mean can be calculated incrementally:

- Let  $X_0, X_1, \dots, X_n$  iid. Then

$$\begin{aligned}\bar{X}_{n+1} &:= \frac{1}{n+1} \sum_{i=1}^{n+1} X_i = \frac{1}{n+1} \left\{ \sum_{i=1}^n X_i + X_{n+1} \right\} \\ &= \frac{n}{n+1} \left\{ \frac{1}{n} \sum_{i=1}^n X_i \right\} + \frac{1}{n+1} X_{n+1}\end{aligned}$$

$$(*) \quad \bar{X}_{n+1} = \left(1 - \frac{1}{n+1}\right) \bar{X}_n + \frac{1}{n+1} X_{n+1}$$

a.s.  $\rightarrow \mathbb{E}X$  as  $n \rightarrow +\infty$ .

More generally, a.s. convergence of  $\bar{X}_n$  to  $\mathbb{E}X$  holds for

$$(**) \quad \bar{X}_{n+1} = (1 - \alpha_n) \bar{X}_n + \alpha_n X_{n+1},$$

provided the sequence  $\{\alpha_n\}$  of positive numbers satisfies

$$\sum_{j \geq 0} \alpha_j = +\infty \quad ; \quad \text{and} \quad \sum_{j \geq 0} \alpha_j^2 < +\infty.$$

(see e.g. thm 14.5 in Mohri & al, Foundations of Machine Learning).

We may rewrite **(\*\*)** slightly differently:

(26)

$$\begin{aligned}\bar{X}_{n+1} &= \bar{X}_n + \alpha_n (X_{n+1} - \bar{X}_n) \\ &= \bar{X}_n + \alpha_n \left( \underbrace{X_{n+1} - \mathbb{E}X}_{\substack{!! \\ Z_{n+1} \\ \text{(zero mean) (iid)}}} + \underbrace{\mathbb{E}X - \bar{X}_n}_{\substack{!! \\ f(\bar{X}_n) \\ \text{with } f(x) = \mathbb{E}X - x}} \right) \\ &= \bar{X}_n + \alpha_n (f(\bar{X}_n) + Z_{n+1}) \quad \nearrow F(x, z) := f(x) + z \\ \bar{X}_{n+1} &= \bar{X}_n + \alpha_n F(\bar{X}_n, Z_{n+1}) \quad \boxed{(***)}\end{aligned}$$

a.s.  $\rightarrow$  to the root of  $f(x) = \mathbb{E}_z \{ F(x, z) \}$

since  $\mathbb{E}_z Z = 0$

&  $f(x) = \mathbb{E}X - x$

&  $F(x, z) := f(x) + z$

Convergence of the sequence  $\{\bar{X}_n\}$  to the root of  $f(x) = \mathbb{E}_z \{ F(x, z) \}$  holds under general assumptions on  $F$ , and on the sequence  $\alpha_n$ .

$\hookrightarrow$  aka. Stochastic Approximation Algorithms

$\hookrightarrow$  dates back to Robbins & Monro (1951).

Ex: We prove in the appendix of the Chapter SL = GRADIENT DESCENT ALGORITHMS that assuming

(a)  $F: \mathbb{R}^2 \rightarrow \mathbb{R}$  s.t.  $\forall x \in \mathbb{R} \quad |F(x, z)| \leq K$   
 $\forall z \in \mathbb{R} \quad (\text{boundedness})$

(b)  $Z_1, \dots, Z_n, \dots$  iid with distribution  $G$ .

(c)  $f(x) = \mathbb{E}_z \{ F(x, z) \}$  is continuous, strictly  $\downarrow$ , with a unique  $x^* \in \mathbb{R}$  such that  $f(x^*) = 0$ .

(d)  $\alpha_n > 0$ ,  $\sum_{j \geq 0} \alpha_j = +\infty$ ,  $\sum_{j \geq 0} \alpha_j^2 < +\infty$

(e)  $X_0$  is square integrable.

Then the sequence  $X_{n+1} = X_n + \alpha_n F(X_n, Z_{n+1})$  converges a.s. to  $x^*$ .

Assumptions (a) - (e) can be weakened considerably. For example, instead of boundedness, we may assume that  $F: \mathbb{R}^d \times \mathbb{R} \rightarrow \mathbb{R}$  is such that  $\forall x \in \mathbb{R}^d \quad \mathbb{E} |F(x, Z)| \leq K(1 + \|x\|^2)$ , for some  $K > 0$

• Application of these ideas to MDP.

Recall that the state-value function  $v_\pi$  satisfies Bellman equation:

$$v_\pi(s) = r(s, \pi(s)) + \gamma \sum_{s'} p(s'|s, \pi(s)) v_\pi(s')$$
$$= \mathbb{E}_{S_{t+1}} \{ r(s, \pi(s)) + \gamma v_\pi(S_{t+1}) \mid S_t = s \}$$

$$\mathbb{E}_{S_{t+1}} \{ \underbrace{r(s, \pi(s)) + \gamma v_\pi(S_{t+1}) - v_\pi(s)}_{F(v_\pi, S_{t+1})} \mid S_t = s \} = 0$$

$\uparrow$  vector of dim  $|S|$       $v_\pi = \begin{pmatrix} | \\ v_\pi(s) \\ | \end{pmatrix} \leftarrow s\text{-th row}$

& we are precisely looking for the root of the function  $\mathbb{E} \{ F(v_\pi, S) \mid S = s \}$ ,  $s \in S$ .  
(function of  $v_\pi$ )

⇒ Consider the sequence of estimates:

$$v_\pi(s_t) \leftarrow v_\pi(s_t) + \alpha \{ r(s_t, \pi(s_t)) + \gamma v_\pi(s_{t+1}) - v_\pi(s_t) \}$$

Updates are "on-line": no need to wait until the end of an episode.     temporal differences of value functions.     ⇒ This procedure is known as TD(0).

Algorithm: TD(0)

- Goal: Estimate  $v_\pi$  of a policy  $\pi$  (unknown environment)
  - Init:  $v_\pi(s) \in \mathbb{R}$  arbitrary except  $v_\pi(\text{terminal states}) = 0$
  - Repeat (for each episode)
    - (i) Initialize  $S$
    - (ii) Repeat for each step of episode:
      - $A \leftarrow$  Action given by  $\pi$  for  $S$
      - Take action  $A$ , observe  $R$  &  $S'$
      - $v_\pi(s) \leftarrow v_\pi(s) + \alpha (R + \gamma v_\pi(S') - v_\pi(s))$
      - $S \leftarrow S'$
- Until  $S$  is a terminal state.

"TD updates a guess towards a guess".

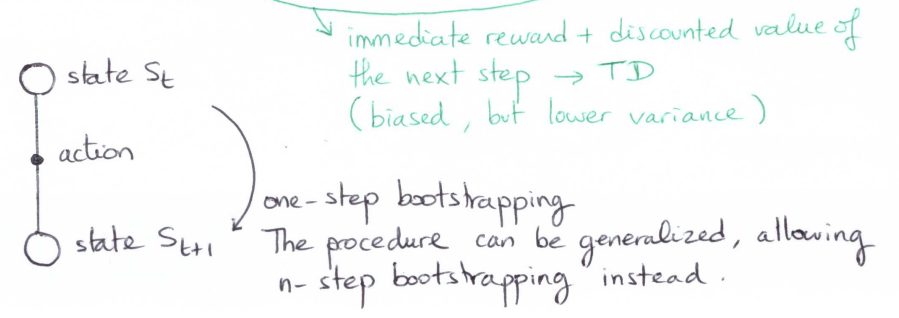
⇒ TD learns at each iteration, before knowing the final outcome. It learns from incomplete episodes by taking advantage of the Markov dynamics of the MDP model.

• Compare:

$$v_\pi(s) = \mathbb{E} ( \underbrace{G_t}_{\text{final reward}} \mid S_t = s )$$

$$= \mathbb{E} ( R_{t+1} + \gamma G_{t+1} \mid S_t = s )$$

$$= \mathbb{E} ( \underbrace{R_{t+1} + \gamma v_\pi(S_{t+1})}_{\text{immediate reward + discounted value of the next step}} \mid S_t = s )$$



n-step bootstrapping simply iterates the previous idea:

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

$$= R_{t+1} + \gamma R_{t+2} + \gamma^2 G_{t+2}, \text{ and so on,}$$

.../...

The n-step return is

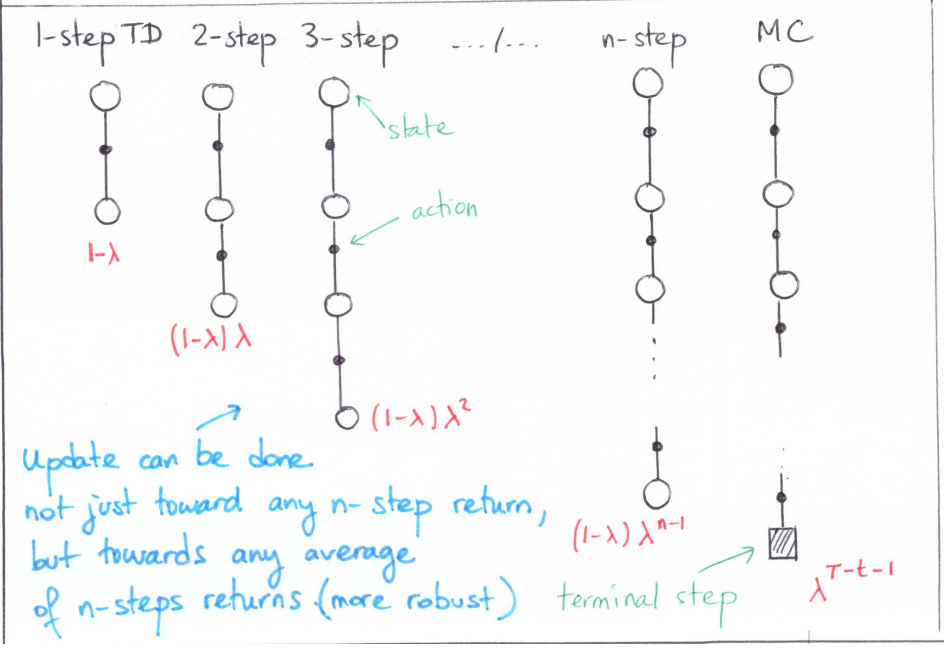
$$G_t^n := R_{t+1} + \gamma R_{t+2} + \dots + \gamma^{n-1} R_{t+n} + \gamma^n v_{\pi}(S_{t+n})$$

= full return truncated after n steps, & corrected for the remaining missing terms by  $v_{\pi}(S_{t+n})$

↳ To compute  $G_t^n$ , we need future rewards  $R_{t+1}, \dots, R_{t+n}$ . The first time these are available is  $t+n \Rightarrow$  the n-step TD algorithm updates according to:

$$v_{\pi}(s_t) \leftarrow v_{\pi}(s_t) + \alpha \{ G_t^n - v_{\pi}(s_t) \}$$

(no updates are made during the first n-1 steps of each episode).

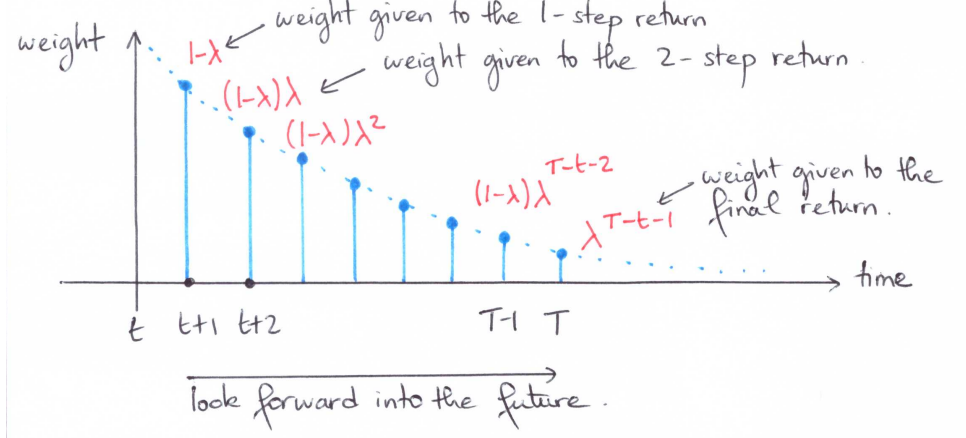


Averaging the n-step returns can be done in many different ways. The  $\lambda$ -return  $G_t^\lambda$  combines all n-step returns  $G_t^n$  using weights  $(1-\lambda)\lambda^{n-1}$ :

$$G_t^\lambda := (1-\lambda) \sum_{n=1}^{\infty} \lambda^{n-1} G_t^n$$

In practice, episodes terminates; and the weight given to the final return is  $\lambda^{T-t-1}$ , to ensure that weights sum to 1.

The  $\lambda$ -return corresponds to a geometrically weighted average.



The TD( $\lambda$ ) algorithm updates the value-function estimate according to

$$v_{\pi}(s) \leftarrow v_{\pi}(s) + \alpha \{ G_t^\lambda - v_{\pi}(s) \}$$

TD( $\lambda$ ) update

Like MC, the TD( $\lambda$ ) update can be made only at the end of an episode.

In practice, take  $\alpha$  small  $\approx 0.1 - 0.3$ .

## IV. MODEL-FREE CONTROL.

31

We consider policy improvement in an MDP environment, with unknown model parameters. We discuss two algorithms

- (i) SARSA
- (ii) Q-learning

### IV.1. SARSA.

- Recall the Bellman equation satisfied by the action-value function  $q_\pi$ :

$$q_\pi(s, a) = r(s, a) + \gamma \sum_{s'} p(s'|s, a) v_\pi(s') \quad (\text{p.8})$$

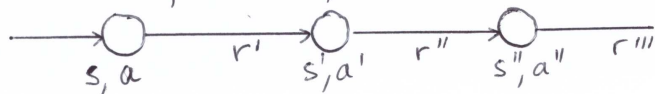
If we knew the model dynamics, we could improve on policy  $\pi$  by acting greedily on  $q_\pi$ ; and select

$$\begin{aligned} \pi'(s) &= \operatorname{argmax}_a q_\pi(s, a) \quad (\text{p.12}) \\ &= \operatorname{argmax}_a \left\{ r(s, a) + \gamma \sum_{s'} p(s'|s, a) v_\pi(s') \right\} \quad (*) \end{aligned}$$

In particular, this ensures that  $v_{\pi'}(s) \geq v_\pi(s)$ . (p.16)

Making this process iterative yields the policy iteration algorithm, of which the value iteration algorithm can be seen as a special case.

- In an unknown environment, the max in (\*) cannot be calculated, and to improve (& evaluate) a policy  $\pi$ , one must learn from experience. Consider a deterministic policy  $\pi$ . Generate trajectories / episodes:



The sequence of observed values (states + rewards) associated with policy  $\pi$  is used to estimate  $q_\pi$  using a stochastic approximation algorithm. Bellman equation for  $q_\pi$  can be rewritten

32

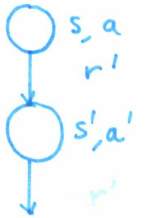
$$\mathbb{E}_{S_{t+1}} \left\{ r(s, a) + \gamma \underbrace{v_\pi(S_{t+1})}_{q_\pi(S_{t+1}, \pi(S_{t+1}))} - q_\pi(s, a) \mid S_t = s, A_t = a \right\} = 0$$

$$\mathbb{E}_{S_{t+1}} \left\{ r(s, a) + \gamma q_\pi(S_{t+1}, \pi(S_{t+1})) - q_\pi(s, a) \mid S_t = s, A_t = a \right\} = 0$$

Consider

$$q(s, a) \leftarrow q(s, a) + \alpha \left\{ r' + \gamma q(s', a') - q(s, a) \right\} \quad (*)$$

where  $\begin{cases} a = \pi(s) \\ a' = \pi(s') \end{cases}$  ← deterministic here, but can be stochastic in general.



Keep iterating until convergence, and then act greedily: select  $\pi'$  such that  $\pi'(s) = \operatorname{argmax}_a q_\pi(s, a)$  ← your estimate

We can improve on this procedure in several directions.

- (a) Speed up the procedure:

No need to wait until converge to act greedily on  $q_\pi$ .

Act greedily after each iteration / update of the action value function:

- Given  $\pi =$  current policy, select  $a \sim \pi$ ,  $a' \sim \pi$
- update  $q(s, a)$  [one iteration of (\*)]
- update policy  $\pi \rightarrow \pi'$  by acting greedily on the new  $q$



(b) Exploration vs Exploitation.

Some pairs (state, action) will not be seen/evaluated by acting greedily  $\Rightarrow$  we are unsure that we end up with the optimal policy  $\Rightarrow$  local minima.

The simplest idea for ensuring a continual exploration is to select the greedy action with probability  $1-\epsilon$ .

$$\pi'(a|s) = \begin{cases} \frac{\epsilon}{|A|} + 1 - \epsilon & \text{if } a = \underset{a'}{\operatorname{argmax}} q(s, a') \\ \frac{\epsilon}{|A|} & \text{otherwise} \end{cases}$$

The  $|A|$  actions are tried with non-zero probability.

$\hookrightarrow$  Terminology: we say that we act " $\epsilon$ -greedily".

Note that for any  $\epsilon$ -greedy policy  $\pi$ , the  $\epsilon$ -greedy policy  $\pi'$  with respect to  $q_\pi$  is an improvement; i.e.  $v_{\pi'}(s) \geq v_\pi(s)$ .

Indeed, first note that

$$\begin{aligned} q_{\pi'}(s, \pi'(s)) &= \sum_{a \in A} \pi'(a|s) q_\pi(s, a) \\ &= \frac{\epsilon}{|A|} \sum_a q_\pi(s, a) + (1-\epsilon) \max_a q_\pi(s, a) \\ &\geq \frac{\epsilon}{|A|} \sum_a q_\pi(s, a) + (1-\epsilon) \sum_a \left\{ \frac{\pi(a|s) - \frac{\epsilon}{|A|}}{1-\epsilon} \right\} \times q_\pi(s, a) \\ &= \sum_a \pi(a|s) q_\pi(s, a) \\ &= v_\pi(s). \end{aligned}$$

The max is  $\geq$  than any weighted sum of its components, provided the weights are  $\geq 0$  & sum to 1 (indeed, weights are  $\geq 0$  since  $\pi$  is an  $\epsilon$ -greedy policy & sum to 1 since  $\sum_a \pi(a|s) = 1$ )

+ use the result on page 16 to conclude that  $v_{\pi'}(s) \geq v_\pi(s)$ . ■

x Algorithm : SARSA.

x Goal : policy improvement in unknown environment

x Init :  $q \leftarrow q_0$  for an arbitrary initial  $q_0$ .

x Repeat (until convergence)  $t = 0, 1, \dots$

$s \leftarrow$  Select state

$a \leftarrow$  Select action  $\sim \pi$  from  $q$ ,  $\epsilon$ -greedily

For each epoch  $t$ , do

$r' \leftarrow$  reward ( $s, a$ )

$s' \leftarrow$  next state ( $s, a$ )

$a' \leftarrow$  select action  $\sim \pi_t$  from  $q$ ,  $\epsilon_t$ -greedily

$$q(s, a) \leftarrow q(s, a) + \alpha_\epsilon (r' + \gamma q(s', a') - q(s, a))$$

$s \leftarrow s'$

$a \leftarrow a'$

x Return  $q$

Remarks: (i) The name of the algorithm comes from the order of successive actions taken  $s - a - r' - s' - a'$ .

(ii) SARSA is an example of an on-policy algorithm; meaning that the update on  $q$  depends on the learning policy (action  $a'$  is selected from  $\pi$ ). In contrast, the Q-learning algorithm presented in the next section is an example of an off-policy algorithm.

(iii) Convergence of SARSA.

SARSA converges to the optimal action-value function

provided

(35)

- $\sum \alpha_t = +\infty$  ;  $\sum \alpha_t^2 < +\infty$  (convergence of the stochastic approximation algorithm  $\rightarrow$  Robbins-Monro)
- Greedy in the Limit with Infinite Exploration (GLIE) sequence of policies  $\pi_t(a|s)$

meaning: we must explore (all pairs  $(s, a)$  must be visited infinitely many times); but we must be greedy in the limit to satisfy Bellman & get to the optimal policy ( $\equiv$  stop exploring).

Take e.g.  $\epsilon_t \sim \frac{1}{t}$ .

(iii) Consider instead n-step returns  $R_{t+1} + \gamma R_{t+2} + \dots + \gamma^n q(S_{t+n}, A_{t+n})$

$\hookrightarrow$  n-step SARSA.

+ consider weighted averages of all n-step return

$\hookrightarrow$  SARSA( $\lambda$ ) algorithm

$\leftarrow$  Same idea as going from TD(0) to TD( $\lambda$ ).

## IV. 2. Q-learning.

The alternative to SARSA is to start from Bellman optimality equation for  $q_*$ :

$$q_*(s, a) = r(s, a) + \gamma \sum_{s'} p(s'|s, a) \max_{a'} \{q_*(s', a')\}$$

$\Leftrightarrow$

$$\mathbb{E}_{S_{t+1}} \left\{ r(s, a) + \gamma \max_{a'} (q_*(S_{t+1}, a')) - q_*(s, a) \mid S_t = s, A_t = a \right\} = 0$$

$\leftarrow$  Solve this recursively using a stochastic approx. algorithm.

Consider

(36)

$$q(s, a) \leftarrow q(s, a) + \alpha \{ r' + \gamma \max_{a'} q(s, a') - q(s, a) \}$$

An "off-policy" algorithm: max is taken over all actions  $a'$ :  $a'$  is not selected from the current policy (compare with SARSA).

x Algorithm: Q-learning

x Goal: policy improvement in an unknown environment

x Init:  $q \leftarrow q_0$  for an arbitrary initial  $q_0$

x For  $t = 0, 1, \dots$ , do

$s \leftarrow$  select state

For each epoch  $t$  do

$a \leftarrow$  select action  $\sim \pi$  from  $q$ ,  $\epsilon_t$ -greedily

$r' \leftarrow$  reward  $(s, a)$

$s' \leftarrow$  next state  $(s, a)$

$$q(s, a) \leftarrow q(s, a) + \alpha_t \{ r' + \gamma \max_{a'} q(s', a') - q(s, a) \}$$

$s \leftarrow s'$

x Return  $q$

$\leftarrow$  To ensure convergence of the Q-learning algorithm, one should make sure that each pair  $(s, a)$  is visited infinitely many times. (+  $\sum \alpha_t = +\infty$ ,  $\sum \alpha_t^2 < +\infty$ ).

## V. SCALING UP

(37)

### V.1. Value Function Approximation.

- So far, given a policy  $\pi$ , a single value of  $v_\pi(s)$  or  $q_\pi(s, a)$  is stored for each  $s$  (or each pair  $(s, a)$ ).

"look-up" tables.

This becomes impractical when the state-space is large.

Ex: Go has  $\approx 10^{170}$  states.

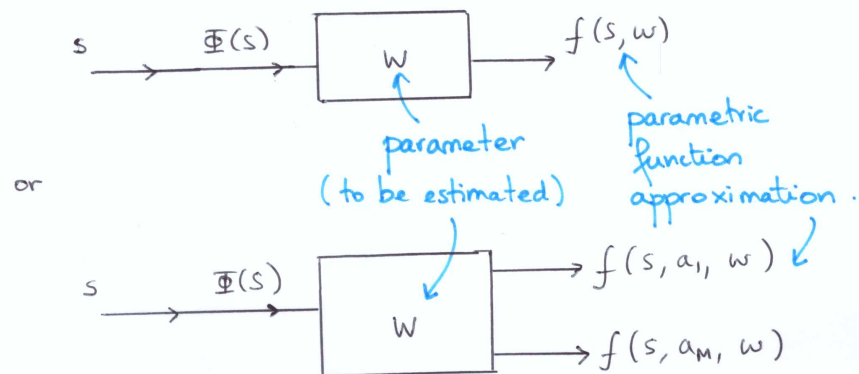
Backgammon  $\approx 10^{20}$  states.

- Idea: Approximate the value function (resp. the action-value function) & generalize from seen states to unseen states

Consider the mapping  $\Phi: \mathcal{S} \rightarrow \mathbb{R}^d$ , with  $d$  relatively small.

a "feature map"

$\forall s \in \mathcal{S}$ ,  $\Phi(s)$  represents some compressed knowledge about  $s$ .  
The low-dimensional feature vector is then used to construct a function  $f(s, w) \equiv$  approximation of  $v_\pi(s)$ .



The problem of estimating a good value of  $w$  learning the input (s) - output ( $v_\pi(s)$ ) relationship can be expressed as a supervised learning problem. (38)

- Ex: (i)  $w$  = vector of parameters in a linear model  
 $f(s, w) = \langle \Phi(s), w \rangle = \Phi(s)^t w$   
Special case:  $\Phi: \mathcal{S} \rightarrow \mathbb{R}^{1 \times |S|}$ , so that  $\Phi$  is the mapping defined by  
 $\Phi(s_j) = \begin{pmatrix} 1 \\ \vdots \\ \mathbf{1} \\ \vdots \\ 1 \end{pmatrix} \leftarrow j$  ;  $\mathcal{S} = \{s_1, \dots, s_{|S|}\}$

Then  $f(s_j, w) = w_j =$  value of state  $s_j$ .  
 $\Rightarrow$  Look-up tables are a special case of approximation methods.

- (ii)  $w$  = set of weights in a neural network.
- (iii) Any other learning algorithm, such as decision trees, nearest-neighbours, ... However, models with differentiable properties are usually preferred, as gradient descent algorithms can be used to estimate  $w$ .

- To start, suppose that we are in presence of an oracle: for each state  $s$ , we know its true label  $v_\pi(s)$ .  
Then given the "training data"  $\{(s_1, v_\pi(s_1)), \dots, (s_T, v_\pi(s_T))\}$ , gathered from experience ( $\equiv$  MC learning, TD learning setting), a (stochastic) gradient method applied to the (empirical) square loss can be used to iteratively update the parameters:

$$w^{(t+1)} = w^{(t)} - \alpha \nabla_w \left\{ \frac{1}{2} (v_\pi(s_t) - f(s_t, w))^2 \right\} \Big|_{w=w^{(t)}} \quad (39)$$

$$= w^{(t)} + \alpha (v_\pi(s_t) - f(s_t, w^{(t)})) \nabla_w f(s_t, w) \Big|_{w=w^{(t)}}$$

For a linear model,  $\nabla_w f(s_t, w) = \Phi(s_t)$ , and the update reduces to:

$$w^{(t+1)} = w^{(t)} + \alpha (v_\pi(s_t) - f(s_t, w^{(t)})) \Phi(s_t).$$

- In practice however, there are no oracle / supervisor. We must substitute a target for  $v_\pi(s)$ , estimated from experience. The most common scenarios are:

→ MC  $v_\pi(s_t) \rightarrow G_t$  (page 4)

→ TD(0)  $v_\pi(s_t) \rightarrow R_{t+1} + \gamma f(s_{t+1}, w^{(t)})$

→ TD( $\lambda$ )  $v_\pi(s_t) \rightarrow G_t^\lambda$  (page 30)

⇒ In the context of MC learning, the training data is  $\{(s_1, G_1), (s_2, G_2), \dots, (s_T, G_T)\}$ , while in the context of TD(0) learning, we consider

$\{(s_1, R_2 + \gamma f(s_2, w)), (s_2, R_3 + \gamma f(s_3, w)), \dots, (s_{T-1}, R_T)\}$

And cycle through the data / use mini-batches (but randomly)

- Control: These updates can be combined with  $\epsilon$ -greedy exploration for policy improvement.

↳ Need to consider approximation for the action-value function:

$$w^{(t+1)} = w^{(t)} + \alpha \left( \overbrace{R_{t+1} + \gamma f(s_{t+1}, A_{t+1}, w^{(t)}) - f(s_t, A_t, w^{(t)})}^{\text{TD(0) update "on-line"}} \right) \times \nabla_w f(s_t, A_t, w) \Big|_{w=w^{(t)}}$$

## II.2. Policy Gradient Methods.

(40)

- An alternative to the action-value method presented in the previous section is to directly parametrize the policy:

$$\pi_\theta(a|s) = \mathbb{P}(A_t = a | S_t = s; \theta)$$

↖ Parameter to be learnt.

↳ Compare with the parametrization of  $v_\pi(s) \rightarrow f(s, w)$   
 $q_\pi(s, a) \rightarrow f(s, a, w)$

The policy is not represented explicitly; but only implicitly once  $q_\pi$  is estimated.

When parametrizing the policy, we still learn a value function to learn the policy parameter  $\theta$ ; but not to select actions.

- value-based → learnt value function / implicit policy
- policy-based → no value function / learnt policy
- actor-critic → learnt value function / learnt policy.

$\theta$  is learnt based on the gradient of some performance measure  $J(\theta)$ , that we wish to maximize:

$$\theta_{t+1} = \theta_t + \alpha \widehat{\nabla J(\theta_t)}$$

gradient ascent

Stochastic estimate of the gradient of  $J$  evaluated at  $\theta_t$ .  
 "On average, points in the right direction".

Such algorithms are known as Policy Gradient Methods (whether or not a value function is learnt)

- The policy is usually parametrized using a soft-max representation =

$$\pi_{\theta}(a|s, \theta) := \frac{e^{h(s,a,\theta)}}{\sum_{a'} e^{h(s,a',\theta)}} \quad 1$$

where  $\rightarrow h(s,a,\theta) = \theta^t \Phi(s,a) =$  linear in features

$\rightarrow h(s,a,\theta) =$  output of a deep neural network  
( $\rightarrow$  Alpha Go).

Consider the episodic case: the performance measure  $J(\theta)$  is the value of the start state of each episode:  $J(\theta) = v_{\pi_{\theta}}^i(s_0)$   
We need to compute the gradient of  $J$ : non-random.

Recall that  $v_{\pi}(s) = \sum_a \pi(a|s) q_{\pi}(s,a)$  (p. 7/8)

$$\begin{aligned} \Rightarrow \nabla v_{\pi}(s) &= \sum_a \left( \nabla \pi(a|s) q_{\pi}(s,a) + \pi(a|s) \nabla q_{\pi}(s,a) \right) \\ &= \sum_a \left( \nabla \pi(a|s) q_{\pi}(s,a) + \pi(a|s) \times \right. \\ &\quad \left. \nabla \left( r(s,a) + \gamma \sum_{s'} p(s'|s,a) v_{\pi}(s') \right) \right) \quad \text{(p. 8)} \\ &= \sum_a \left( \nabla \pi(a|s) q_{\pi}(s,a) \right. \\ &\quad \left. + \pi(a|s) \sum_{s'} p(s'|s,a) \nabla v_{\pi}(s') \right) \\ &= \sum_a \left( \nabla \pi(a|s) q_{\pi}(s,a) \right. \\ &\quad \left. + \pi(a|s) \sum_{s'} p(s'|s,a) \left\{ \sum_{a'} \nabla \pi(a'|s') q_{\pi}(s',a') \right. \right. \\ &\quad \left. \left. + \pi(a'|s') \sum_{s''} p(s''|s',a') \nabla v_{\pi}(s'') \right\} \right) \end{aligned}$$

$$= \sum_{s'} \left( \underbrace{\sum_a \pi(a|s) p(s'|s,a)}_{=: p(s \rightarrow s', 1, \pi)} \right) \left( \sum_{a'} \nabla \pi(a'|s') q_{\pi}(s',a') \right)$$

$=: p(s \rightarrow s', 1, \pi) =$  probab of going from state  $s$  to state  $s'$  in one step, following policy  $\pi$

Keep unrolling:

terms  $\sum_{s'} p(s \rightarrow s', k, \pi) \sum_{a'} \nabla \pi(a'|s') q_{\pi}(s',a')$  appear.

$$\Rightarrow \nabla v_{\pi}(s) = \sum_{s'} \sum_{k \geq 0} p(s \rightarrow s', k, \pi) \sum_{a'} \nabla \pi(a'|s') q_{\pi}(s',a')$$

Thus

$$\nabla J(\theta) = \sum_s \left( \underbrace{\sum_{k \geq 0} p(s_0 \rightarrow s, k, \pi)}_{=: \eta(s)} \right) \sum_a \nabla \pi(a|s) q_{\pi}(s,a)$$

$\uparrow$  initial deterministic state

$$\begin{aligned} &= \sum_s \eta(s) \sum_a \nabla \pi(a|s) q_{\pi}(s,a) \\ &= \sum_{s'} \eta(s') \sum_s \left( \frac{\eta(s)}{\sum_{s'} \eta(s')} \right) \left( \sum_a \nabla \pi(a|s) q_{\pi}(s,a) \right) \end{aligned}$$

$=: \mu(s)$

fraction of time spent in each state  $s$ .  
 $\mu$  sums to one (aka the "on-policy distribution")

$$\begin{aligned} &= \sum_{s'} \eta(s') \sum_s \mu(s) \left( \sum_a \nabla \pi(a|s) q_{\pi}(s,a) \right) \\ &\propto \sum_s \mu(s) \left( \sum_a \nabla \pi(a|s) q_{\pi}(s,a) \right) \end{aligned}$$

(proportional to) (enough for us  $\rightarrow$  will be absorbed in the learning rate later)

$$\nabla J(\theta) \propto \mathbb{E}_{\pi} \left( \sum_a q_{\pi}(S_t, a) \nabla \pi(a|S_t, \theta) \right)$$

aka the Policy Gradient Theorem.

replacing  $s$  with  $S_t$  will allow us to make use of samples collected from experience to compute  $\nabla J(\theta) \Rightarrow$  we want to replace  $a$  with  $A_t$  as well.

Next, notice that

$$\begin{aligned} \nabla \pi(a|s, \theta) &= \pi(a|s, \theta) \frac{\nabla \pi(a|s, \theta)}{\pi(a|s, \theta)} \\ &= \pi(a|s, \theta) \nabla \ln \pi(a|s, \theta) \end{aligned}$$

Thus

$$\begin{aligned} \nabla J(\theta) &\propto \mathbb{E}_{\pi} \left( \sum_a \pi(a|S_t, \theta) q_{\pi}(S_t, a) \nabla \ln \pi(a|S_t, \theta) \right) \\ &= \mathbb{E}_{\pi} \left( q_{\pi}(S_t, A_t) \nabla \ln \pi(A_t|S_t, \theta) \right) \\ &= \mathbb{E}_{\pi} \left( \underbrace{G_t}_{\text{return}} \nabla \ln \pi(A_t|S_t, \theta) \right) \end{aligned}$$

↑  
return =

We can make use of this expression using a SG algorithm.

⇒ Updates are  $\theta_{t+1} = \theta_t + \alpha G_t \nabla \ln \pi(A_t|S_t, \theta_t)$ .

MC policy - gradient control.

(generate an episode  $s_0, A_0, R_1, S_1, A_1, \dots, S_{T-1}, A_{T-1}, R_T$  following  $\pi(\cdot|\cdot, \theta)$ , and for each step of the episode  $t=0, \dots, T-1$ , update  $\theta \leftarrow \theta + \alpha G \nabla \ln \pi(A_t|S_t, \theta)$ )

& repeat for each episode.

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

(taking  $\gamma=1$ )

↑

This approach can be improved, see. e.g. Actor-Critic Methods.