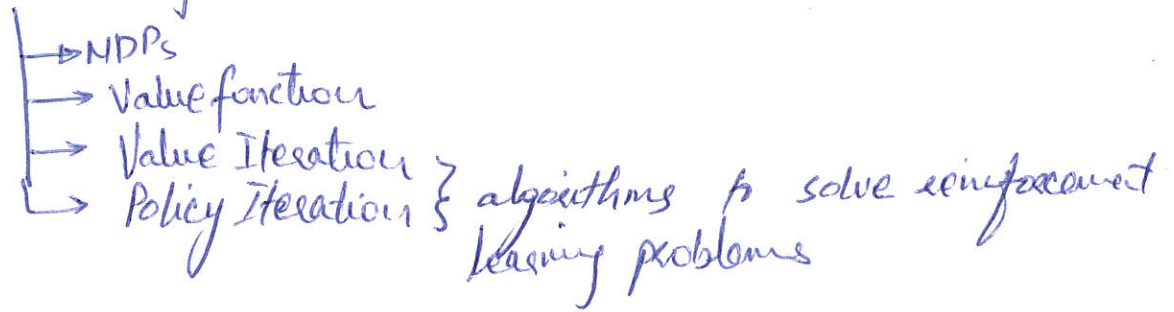


# Reinforcement learning (4/16)

①

## \* Reinforcement learning



## \* CREDIT Assignment Problem

Suppose you are playing a chess. You lost a game at 60<sup>th</sup> move. CAP tries to associate reward with every move you take. However, you get to know about the outcome of chess after 60<sup>th</sup> move.

## \* Formalizing RL

RL model the world using Markov Decision Process (MDPs)

## MDP

is a five tuple:  $(S, A, \{P_{sa}\}, \gamma, R)$

$S$ : set of states (in helicopter example it will be the possible orientation helicopter could be in)

$A$ : set of action (set of possible controller configurations)

$P_{sa}$ : State-transition distribution

$$\sum_{s'} P_{sa}(s') = 1, \quad P_{sa}(s') \geq 0$$

( $P_{sa}(s')$  is the probability of transitioning in state  $s'$  if we take action "a" in state  $s$ ).

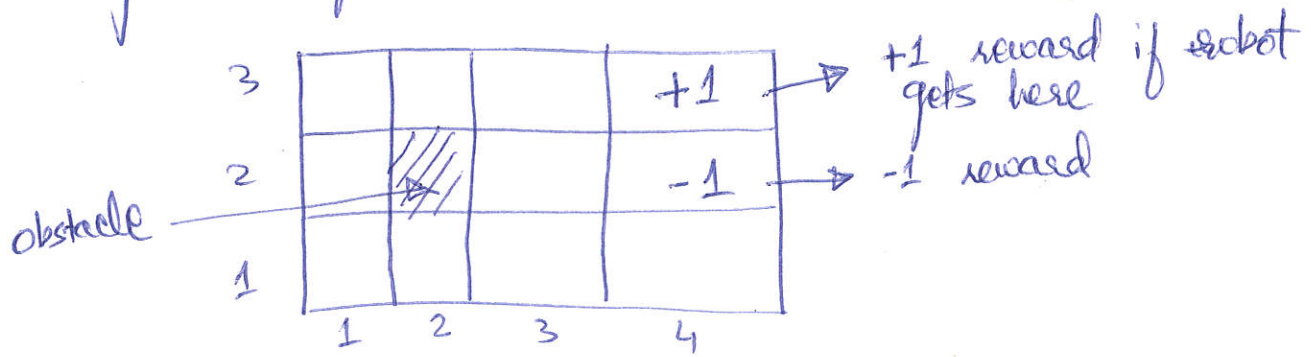
$\gamma$ : Discount factor  $0 \leq \gamma < 1$

$R$ : reward function  $R: S \rightarrow \mathbb{R}$  (real-numbers)

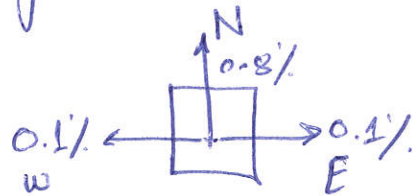


# Running Example for Peter Russell book (Robot Navigation Example)

(3)



- \* We have 11 states here (Robot can be in any of these places)
- \* Robot can take 4 actions:  $A = \{N, S, E, W\}$
- \* Robot is <sup>stochastic</sup> dynamic as well, so we will model it as well



(If we command the robot to move "N", then there is 80% chance that robot will move N & 10% that it will deviate toward left & right)

\* Let's write state transition probability when robot is at location (3,1) & we ask the robot to move to (3,2)

$$P_{(3,1), \text{move north}}((3,2)) = 0.8$$

$$P_{(3,1), N}((4,1)) = 0.1$$

$$P_{(3,1), N}((2,1)) = 0.1$$

$$P_{(3,1), N}((3,3)) = 0$$

etc

$$P_{(3,1), N}^S((4,1)) = 0.1$$

If robot is at state (3,1) and takes action "N", then the probability that robot will transition into state (4,1) is 0.1 or 10%.



\* Reward (R)

(5)

We have  $R(4,3) = +1$

$R(4,2) = -1$

and  $R(s) = -0.02$  (~~that~~ small negative reward for all other states)

Assumption

World ends when robot reaches  $R(4,3)$ . (Zero Cost absorbing state)

How MDPs work

At state  $s_0$

Choose  $a_0$

Get to  $s_1$  by randomly drawing  $s_1: s_1 \sim P_{s_0 a_0}$

Choose  $a_1$

Get to  $s_2 \sim P_{s_1 a_1}$

$\vdots$

$s_0 \xrightarrow{a_0} s_1 \sim P_{s_0 a_0} \xrightarrow{a_1} s_2 \sim P_{s_1 a_1} \rightarrow \dots \text{(end)}$

After a while, robot would have visited sequence of state

$s_0 \quad s_1 \quad s_2 \quad \dots$

And in order to evaluate how well ~~we~~ robot did, we will apply reward function to states & sum them

$R(s_0) + R(s_1) + R(s_2) \dots$

We will sum them up with discount factor, giving total payoff

Total Payoff =  $R(s_0) + \gamma R(s_1) + \gamma^2 R(s_2) + \dots$

$0 \leq \gamma < 1$  (i.e. reward of future state have less weight than their past)





The goal of learning algorithm is  
 to choose actions over time  $(a_0, a_1, \dots)$   
 to maximize the expected value of total payoff:  

$$E[R(s) + \gamma R(s_1) + \gamma^2 R(s_2) + \dots]$$

More Concretely, RL will compute a policy  $\pi$ .  
 Policy  $\pi: S \rightarrow A$   
 $\uparrow \quad \quad \uparrow$   
 state action

An Example of Policy  $\pi$   
Optimal Policy for Robot Navigation

3	→	→	→	+1
2	↑	////	↑	-1
1	↑	←	←	←
	1	2	3	4

→ If you are in state  $(4,1)$ , then take "left" action  
 i.e.  $\pi((4,1)) = W$   
 $A = \{W, N, E, S\}$

\* This policy will maximize the expected value of total payoff

Let's define the following:  
 $V^x, V^*, \pi^*$

$\pi^*$  seems to be optimal policy, but it is not.  
 It will be optimal as a consequence of our definition.

For any policy  $x$ :

For any  $x$ , define value function  $V^\pi: S \rightarrow \mathbb{R}$





(9)

st  $V^\pi(s)$  is expected total payoff starting in state  $s$ , and execute  $\pi$ .

$$V^\pi(s) = E \left[ R(s_0) + \gamma R(s_1) + \dots \mid \pi, s_0 = s \right]$$

This is sloppy notation as  $\pi$  is not a random variable and we shouldn't be conditioning on it.

As a concrete example, let's consider (bad) policy  $\pi$ , and its value function  $V^\pi$

$\pi$

3	→	→	→	+1
2	↓	///	→	-1
1	→	→	↑	↑
	1	2	3	4

$V^\pi$

3	0.52	0.73	0.77	+1
2	-0.9	///	-0.82	-1
1	-0.88	-0.87	-0.85	-1
	1	2	3	4

$$V^\pi(s) = E \left[ \underbrace{R(s_0)}_{\text{Immediate reward}} + \underbrace{\gamma (R(s_1) + \gamma R(s_2) + \dots)}_{\substack{\text{future rewards} \\ V^\pi(s_1)}} \mid \pi, s = s_0 \right]$$

$$\begin{array}{l} s_0 \rightarrow s \\ s_1 \rightarrow s' \end{array}$$

$$V^\pi(s) = R(s) + \gamma \sum_{s'} \underbrace{P(s'|s, \pi)}_{P_\pi(s')} V^\pi(s') \quad \left\{ \begin{array}{l} \text{we are executing} \\ \text{policy } \pi \end{array} \right.$$

→ Bellman's Equation



$$V^x(s) = R(s) + \gamma \dots V^x(s')$$

$\uparrow$   
 $s'$  is random

$P_{sa}(s')$  where  $a = x(s)$

$P_{sx(s)}(s')$

Bellman's Equation allows us to solve  $V^x(s)$ . Given a policy  $x$ , how to compute  $V^x(s)$ . Answer is Bellman's equation.

Let's make Bellman's Eq more concrete. Let's consider (3,1) state:  
 & let's assume we have a specific policy  $x$  s.t.  $x((3,1)) = N$

$$V^x((3,1)) = R((3,1)) +$$

3				+1
2		///		-1
1			○	
	1	2	3	4

$$\gamma \left[ 0.8 V^x((3,2)) + 0.1 V^x((2,1)) + 0.1 V^x((4,1)) \right]$$

Bellman's equation for state (3,1) with policy  $x$

Similarly, for all 11-states, I can write  $V^x(s)$ .

i.e.  $V^x((1,1)) = \dots$

$V^x((2,1)) = \dots$

$V^x((4,3)) = \dots$

This means we will have 11-variable,  $[V^x((1,1)) \dots V^x((4,3))]$  and 11-constraints. of linear system.



Optimal value function :

(13)

$$V^*(s) = \max_{\pi} V^{\pi}(s)$$

Bellman's Eq/s: (Bellman's eq for best value function)

$$V^*(s) = R(s) + \max_a \gamma \sum_{s'} P_{sa}(s') V^*(s')$$

$$\pi^*(s) = \operatorname{argmax}_a \sum_{s'} P_{sa}(s') V^*(s')$$

The consequence of above equations will be that  $\pi^*$  is the optimal policy.

### VALUE Iteration

This is an algorithm to compute optimal policy using Bellman's Equation.

Initialize  $V(s) = 0 \quad \forall s$

repeat until convergence { For every  $s$ , update  $V(s) := R(s) + \max_a \gamma \sum_{s'} P_{sa}(s') V(s')$

This will make  $V(s) \mapsto V^*(s)$

There could be two methods of updating  $V(s)$

i) Synchronous Update:

$$V(s) := R(s) + \max_a \gamma \sum_{s'} P_{sa}(s') V(s')$$

CALC RHS for all states

then  $V := B(V)$

ii) Asynchronous update:

Use the new value of  $V(s)$  in update equation





# VALUE Iteration example

(15)

3	$s_8$	$s_9$	$s_{10}$	$+   s_{11}$
2	$s_5$	///	$s_6$	$-   s_7$
1	$s_1$	$s_2$	$s_3$	$s_4$
	1	2	3	4

Init  $V(s) = 0$

- $V(1,1) = 0.0$
- $V(2,1) = 0$
- $V(3,1) = 0$
- $V(4,1) = 0$
- $V(1,2) = 0$
- $V(3,2) = 0$
- $V(1,3) = 0$
- $V(2,3) = 0$
- $V(3,3) = 0$
- $V(4,3) = 0$

OR  $V(s_1) = V(s_2) = \dots = V(s_{11}) = 0.0$

## SYNCHRONOUS UPDATE

$$\begin{bmatrix} V(s_1) \\ V(s_2) \\ \vdots \\ V(s_{11}) \end{bmatrix} = \begin{bmatrix} R(s_1) + \max_a \gamma [P_{s_1 a}(s_5) V(s_5) + \dots] \\ \vdots \\ R(s_{11}) + \max_a \gamma [P_{s_{11} a}(s_{11}) V(s_{11}) + \dots] \end{bmatrix}$$

## ASync Update

$$V(s_1) = R(s_1) + \max_a \gamma [P_{s_1 a}(s_5) V(s_5) + P_{s_1 a}(s_2) V(s_2)]$$

$$V(s_2) = R(s_2) + \max_a \gamma [P_{s_2 a}(s_3) V(s_3) + P_{s_2 a}(s_1) V(s_1)]$$

$$V(s_3) = \vdots$$

$$V(s_{11}) =$$



To use this algorithm, let's use  $\gamma = 0.99$

(17)

3				+1
2		///		-1
1				
	1	2	3	4

Let's run value iteration on MDP above. The numbers we get for  $V^*$  are as follows.

$V^* =$

3	0.88	0.90	0.93	+1
2	0.92	///	0.69	-1
1	0.78	0.75	0.71	0.49
	1	2	3	4

If you plug this  $V^*$  into the formula

$$\pi^*(s) = \underset{a}{\operatorname{argmax}} \sum_{s'} P_{sa}(s') V^*(s')$$

$\pi^* =$

3	→	→	→	+1
2	↑	///	↑	-1
1	←	←	←	←
	1	2	3	4

To summarize, we run value iteration to get  $V^*$  (as shown above: table of numbers) and then I compute optimal policy using  $\pi^*(s)$  formula.

In  $\pi^*$ , let's take a look at state (3,1). Why do we select action "W"?  $A: \{W, N, E, S\}$ .

At (3,1)

Moving west:  $\sum_{s'} P_{sa}(s') V^*(s')$

$$= 0.8 \times 0.75 + 0.1 \times 0.69 + 0.1 \times 71 = 0.740$$

Moving North:  $\sum_{s'} P_{sa}(s') V^*(s')$

$$= 0.8 \times 0.69 + 0.1 \times 0.75 + 0.1 \times 0.49 = 0.676$$

(58:40)  
L#15



If Robot is at (3,1), then the expected payoff of moving west (0.74) is higher than moving north (0.676). That's why  $\boxed{X((3,1)) = W}$

At (3,1):

$$\begin{aligned} \text{Moving West: } & \sum_{s'} P_{sa}(s') V^*(s') \\ \text{state: } (2,1) &= (0.8 \times 0.75) + (0.1 \times 0.69) \\ &+ (0.1 \times 0.71) \\ \text{state: } (3,1) & \end{aligned}$$

$$\begin{aligned} \text{Moving North (3,2)} &= \sum_{s'} P_{sa}(s') V^*(s') \\ &= (0.8 \times 0.69) + (0.1 \times 0.75) + (0.1 \times 0.49) \\ &\quad \text{state: } (3,2) \quad \text{state: } (2,1) \quad \text{state: } (4,1) \end{aligned}$$

Another algorithm to find optimal policy for MDP is "Policy Iteration".

## Policy Iteration

Initialize  $\pi$  randomly

Repeat:

- 1) Let  $V_0 = V^\pi$  (Solve Bellman's Eq) (solve system of linear equations)  
if unknown all constants
- 2) Let  $\pi(s) := \arg\max_a \sum_{s'} P_{sa}(s') V(s')$

$$V \mapsto V^* \quad \pi \mapsto \pi^*$$





Step(1) in Policy iteration is expensive. If you have  $N$  state, then you have to solve a <sup>linear</sup> system of  $N$  ~~unknown~~ variables <sup>with</sup>  $N$  constraints.

If your MDP is very large ( $\approx 10$  million states), then it's better to use value iteration than policy iteration as policy iteration would need to solve system of around  $\approx 10$  million equations.

\* Now you can use policy iteration or value iteration to <sup>compute</sup> ~~solve~~ optimal policy for a given MDP.

Now, what if you don't know state-transition probabilities ( $P_{sa}$ )?  
 $\rightarrow$  MDP:  $S, A, P_{sa}, \gamma, R$  (we will estimate  $P_{sa}(s')$  from data)

$$P_{sa}(s') = \frac{\# \text{ times took action "a" in } s, \text{ got to } s'}{\# \text{ times took action "a" in } s}$$

$$P_{sa}(s') = \begin{cases} \frac{1}{|S|} & \text{if } \frac{0}{0} \end{cases}$$

Putting it together

Repeat: {

\* Take actions using  $\pi$  to get experience in MDP

\* Update estimates of  $P_{sa}$

\* Solve Bellman's eq' using VI to get  $V$

\* Update  $\pi(s) = \arg\max_a \sum_{s'} P_{sa}(s') V(s')$

}

} Solving MDPs for those problems where  $P_{sa}$  is not available

