

## Reinforcement Learning (RL)

144

• RL differs from supervised learning in two crucial respects:

1) SL is passive: model observes world and computes "predictions"

RL is active: observes, makes predictions, and executes actions that influence the environment (in desired ways)

2) SL is strongly supervised: for each training instance, the true answer is known

RL is weakly supervised: a feedback signal is often only available after many actions

example: learning chess: SL  $\Rightarrow$  student has a coach who explains after every move whether it was good or bad and why  
RL  $\Rightarrow$  feed back only game outcome "win", "loss", "draw"  $\Rightarrow$  1.5 bits after  $\approx 40$  actions

• feedback in RL is formally defined by "rewards": transition from time step

$t \rightarrow t+1$  produces reward  $R_{t+1}$ :  $R_{t+1} > 0 \hat{=}$  good, gain

$R_{t+1} < 0 \hat{=}$  bad, penalty

$R_{t+1} = 0 \hat{=}$  neutral, no feedback in this time step

- extreme case: executes  $T$  time steps (chess  $T \approx 40$ )

and then gets a single reward, all rewards for  $t < T$  are neutral

$\Rightarrow$  makes learning difficult  $\Rightarrow$  reward shaping  $\hat{=}$  manually introduces auxiliary rewards during training

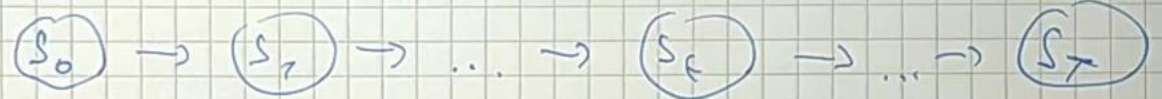


Sub: reward shaping must be used with care, because auxiliary rewards are not present at test time and may cause training to converge to a bad local ~~max~~ optimum  $\Rightarrow$  algorithm "exploits" the shortcomings of the auxiliary rewards, e.g. infinite loops that do not help winning

formal definitions: agent acts in an environment  $\mathcal{S}$

- environment is characterized by a state  $s_t$  at time  $t$ , states evolve in discrete time steps  $s_t \rightarrow s_{t+1}$

$\Rightarrow$  states form a chain



- some environments always reach a terminal state after finite time  $T < \infty$ , e.g. most games and with win or loss  $\hat{=}$  different runs: "episodes"
- others can run indefinitely, e.g. typical in the real world
- transition probabilities of the environment: likelihood for the next state, when we are in a given state

$$p(s_{t+1} = s' \mid s_t = s)$$

$\uparrow$  next chess position       $\uparrow$  current chess position

- crucial assumption: Markov property: the future only depends on the present, not on the past (chess: does not matter how one ends up in position  $s$ , the best next move is always independent of history)

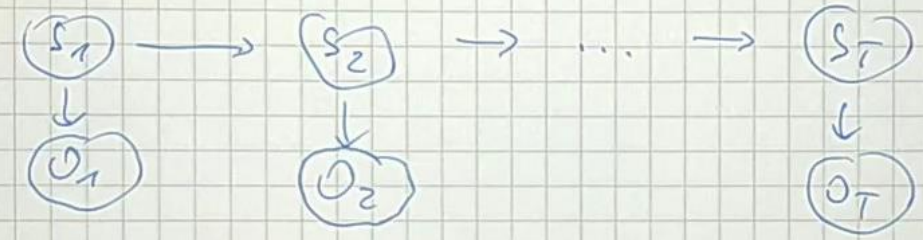
$$s_{t-1} \perp s_{t+1} \mid s_t \Rightarrow \boxed{p(s_1, \dots, s_T) = p(s_1) \prod_{t=2}^T p(s_t \mid s_{t-1})}$$

states form "Markov chain"



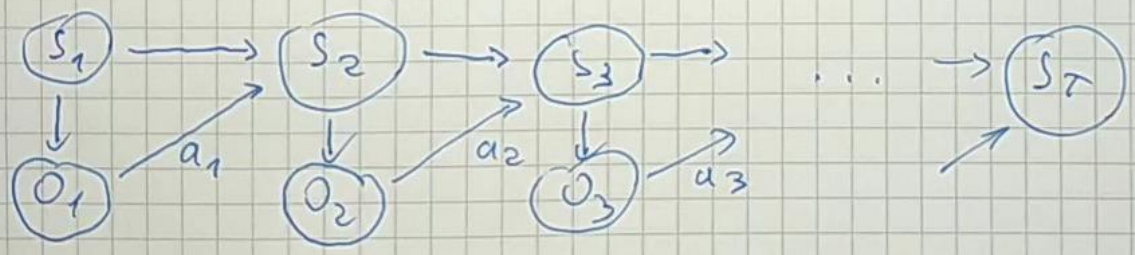
- observations: the agent may not be able to know the present state exactly, but only gets indirect information via observations  $O_t$
- fully observable:  $O_t = S_t$  (example: chess)
- partially observable:  $O_t \subset S_t$  (example: card game)
- state is hidden: observation probability  $p(O_t | S_t)$ 
  - $O_t$  only depends on the present  $S_t$  not the past  $S_{t-1}, \dots$

→ hidden Markov model HMM



examples: speech recognition  
 Poker: reactions of opponents

- so far, we only considered the environment passively, but actually the agent's actions can influence the state transitions



redifine state transitions:  
 $p(S_{t+1}=s' | S_t=s, a_t=a)$

actions are chosen according to the current observations (history plays no role)  $\hat{=}$  policy of the agent  $\pi(O_t)$

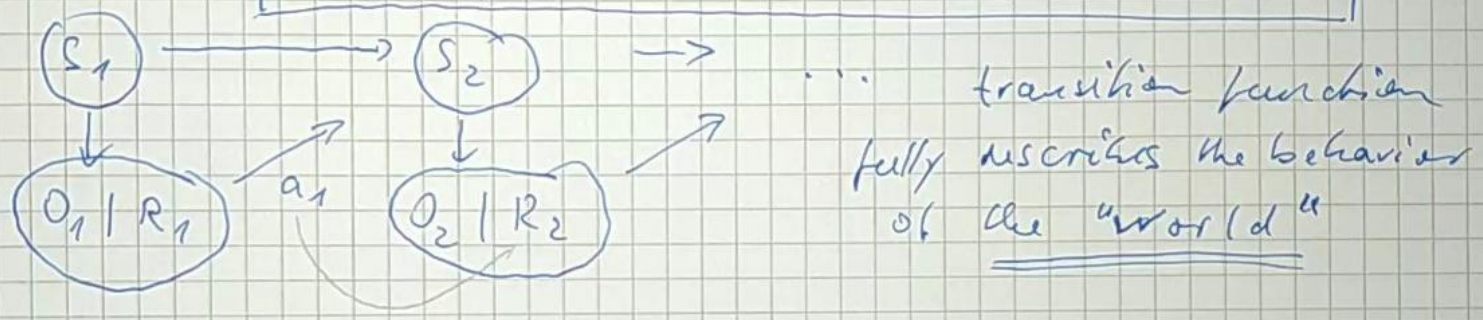
- deterministic policy:  $a_t = \pi(O_t)$  is always the same in state  $O_t = O$
- probabilistic policy:  $a_t \sim \pi(a | O_t)$   $\Rightarrow$  can choose different actions in same situation (good ones with high probability)



- introduce rewards into the model, again ~~deterministic~~ probabilistically, e.g. in general there is no guarantee to receive a specific reward in a given situation, rewards are chosen from a distribution.

convention: reward of transition  $t \rightarrow t+1$  is called  $R_{t+1}$

full transition model: 
$$p(S_{t+1}=s', R_{t+1}=r | S_t=s, a_t=a)$$



the agent is fully described by the observations that receives and the policy to choose actions:  $p(O_t | S_t), \pi(a_t | O_t)$

- simplification in the lecture: consider only fully observable worlds:  $O_t \equiv S_t$   
 $\Rightarrow p(O_t | S_t) = \delta(O_t - S_t)$ ,  $\pi(a_t | S_t)$  conditioned on current state  
 $\Rightarrow$  Markov Decision Process MDP (generalization: POMDP "partially observable MDP:  $O_t \neq S_t$ ")

- Goal of RL: given training data, find the policy  $\pi^*(a_t | s_t)$  that maximizes the total reward

$\pi^*(a_t | s_t)$  unknown optimal policy,  $\hat{\pi}(a_t | s_t) \approx \pi^*(a_t | s_t)$



- when the state transitions and/or policy are non-deterministic, the same initial state may result in many different games with different rewards  
 $\Rightarrow$  take the expected reward in the optimization  $\hat{=}$  value function

$$V_{\pi}(s) = \mathbb{E}_{p, \pi} \left[ \sum_{t'=t+1}^T R_{t'} \mid s_t = s \right]$$

$\uparrow$   
 start state  
 at time  $t$

when  $T$  is very big or  $T \rightarrow \infty$ ,  $\sum R_{t'}$  may diverge  $\Rightarrow$   
 discounted value function with discount factor  $\gamma$  (hyperparameter)

$$V_{\pi}(s) = \mathbb{E}_{p, \pi} \left[ \sum_{t'=t+1}^{\infty} \gamma^{t'-t-1} R_{t'} \mid s_t = s \right]$$

$0 < \gamma \leq 1$  :  $\gamma^{t'-t-1}$   $\gamma^0, \gamma^1 < 1, \dots, \gamma^k \ll 1 \quad k \rightarrow \infty$   
 $\Rightarrow$  future rewards are downweighted

horizon for reward calculation: narrow when  $\gamma$  small

- optimal policy maximizes the expected reward

$$\left[ \begin{aligned} \pi^*(s) &= \max_{\pi} \arg \max_{\pi} V_{\pi}(s) \\ V^*(s) &= V_{\pi^*}(s) \end{aligned} \right] \quad (\text{deterministic policy case})$$

- simplification: introduce "return" : sum of rewards :

$$G_t = \sum_{t'=t+1}^{\infty} R_{t'}$$

or

$$G_t = \sum_{t'=t+1}^{\infty} \gamma^{t'-t-1} R_{t'}$$

$\xrightarrow{\text{equal}}$

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



- rewrite value function with "return":

$$V_{\pi}(s) = \mathbb{E}_{P, \pi} [G_t \mid S_t = s] = \mathbb{E}_{P, \pi} [R_{t+1} + \gamma G_{t+1} \mid S_t = s]$$

- assume that state and action spaces are discrete and finite

⇒ replace expectation with a sum

$$V_{\pi}(s) = \sum_a \pi(a|s) \sum_{s', r} p(s', r | s, a) \left[ r + \underbrace{\mathbb{E}_{P, \pi} [V_{\pi}(G_{t+1}) \mid S_{t+1} = s']}_{\substack{\text{immediate reward} \\ \text{expected return of the rest of the game, starting in state } s'}} \right]$$

$a$  ← all permissible actions in state  $s$      
  $s', r$  ← all possible rewards from transition  $s \rightarrow s'$

now  $\mathbb{E}_{P, \pi} [\gamma G_{t+1} \mid S_{t+1} = s'] = V_{\pi}(s') \cdot \gamma$

⇒ define  $V_{\pi}(s)$  recursively in terms of  $V_{\pi}(s')$

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

Bellman equation for value function  $\hat{=}$  enforces self-consistency of the value function

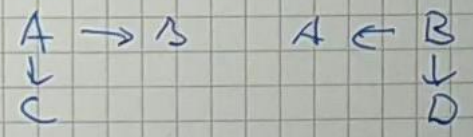
since states are discrete, we have finitely many numbers  $V_{\pi}(s)$

for a given policy  $\pi \Rightarrow$  Bellman eq. is just a linear system

example simple game with 4 squares, allowed actions

A = start	B
C = monster	D = treasure

game terminates when in state C or D



deterministic rewards:  $V_{A \rightarrow C} = -100$ ,  $V_{B \rightarrow D} = 100$

$V_{A \rightarrow B} = V_{B \rightarrow A} = -1$  (avoid infinite loop  $A \leftrightarrow B$ )



- define the "world"  $p(s', v | s, a)$ :

- easy to see optimal policy:

$$\pi^*(s=A) = \rightarrow$$

$$\pi^*(s=B) = \downarrow$$

[  $\pi^*(s=C), \pi^*(s=D)$  undefined,  
because game over ]

- write out Bellman equations for  $\pi^*$   
( $\gamma = 1$ ) no discount)

$$v^*(s=A) = 0.8 [-1 + v^*(s=B)] + 0.2 [-100 + \underbrace{v^*(s=C)}_{=0, \text{ because game over}}]$$

$$v^*(s=B) = 0.8 [100 + \underbrace{v^*(s=D)}_{=0}] + 0.2 [-1 + v^*(s=A)]$$

$$20.8 = -v^*(s=A) + 0.8 v^*(s=B)$$

$$-79.8 = 0.2 v^*(s=A) - v^*(s=B)$$

$$v^*(s=A) \approx 51$$

$$v^*(s=B) \approx 90$$

- easily calculate that any other policy has lower value, for example

$$\pi^0(s=A) = \downarrow, \pi^0(s=B) = \leftarrow$$

$$v^0(s=A) = -91, v^0(s=B) = -54$$

$s'$	$r$	$s$	$a$	$p(\dots)$
B	-1	A	$\rightarrow$	0.8
C	-100	A	$\rightarrow$	0.2
(action ignored)				
C	-100	<del>A</del>	$\downarrow$	0.8
B	-1	A	<del><math>\rightarrow</math></del>	0.2
D	100	<del>A</del>	$\downarrow$	0.8
A	-1	B	$\downarrow$	0.2
A	-1	B	$\leftarrow$	0.8
D	100	B	$\leftarrow$	0.2