

Reinforcement Learning

Prof. Christian Bauckhage



lecture 07

foundations of reinforcement learning

given all our previous lectures, we can finally study reinforcement learning

in this lecture, we therefore have our first look at the modeling tools which are most commonly used to formalize the reinforcement learning problem

in particular, we introduce *Markov decision processes*, *(expected) returns*, *stochastic policies*, *value functions*, and *quality functions*

in our upcoming lectures, we will study all these concepts in great detail



advise:

work through these slides! they lay the foundation for everything to come!

outline

recap

reinforcement learning

- Markov decision processes

- policies and value functions

- Bellman equations

summary

recap

discrete state space

⇔ a finite or countably infinite set $\mathcal{S} = \{s_1, s_2, \dots\}$

stochastic process

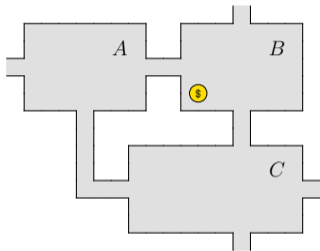
⇔ a set of rvs $\{S_t \mid t \in \mathbb{N}\}$ where each $S_t \in \mathcal{S}$

discrete time Markov chain (DTMC)

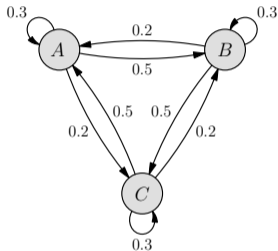
⇔ a stochastic process such that for all $t \in \mathbb{N}$

$$p(S_t \mid S_{t-1}, \dots, S_0) = p(S_t \mid S_{t-1}) = p(S_t = s' \mid S_{t-1} = s) \equiv p(s' \mid s)$$

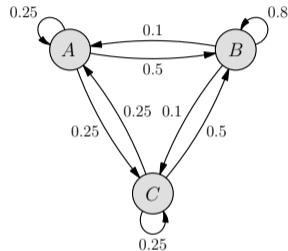
Markov chains provide simple (“first-order”) models of behavior



an environment



λ_1 : patrolling the map



λ_2 : guarding the treasure

reinforcement learning

reinforcement learning

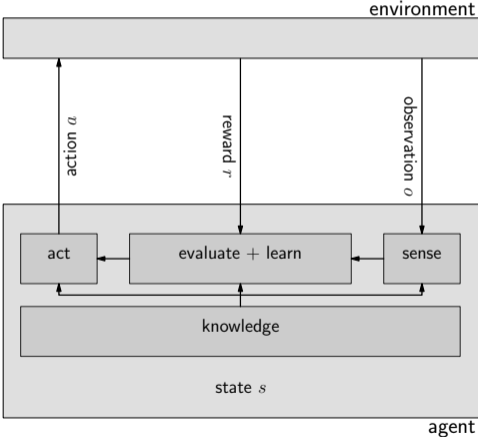
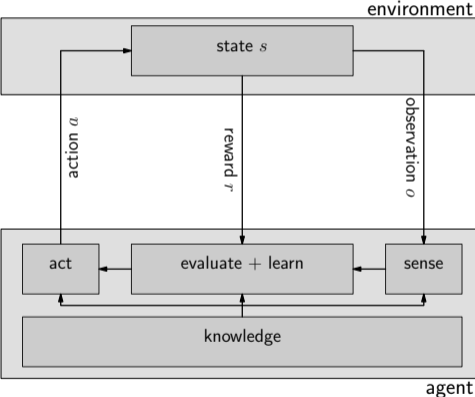
- ⇔ the problem faced by an *agent* that uses *informed trial-and-error explorations* to learn how to *behave* in a dynamic environment
- ⇔ the problem faced by an agent that uses trial-and-error to learn to decide *what* to do *when* in order to achieve a *long-term goal*

there are two major strategies for computational behavior learning

- 1) *search* the space of all behaviors to find those that work well (e.g. using MCTS, genetic algorithms, . . .)
- 2) **apply *statistical techniques* and *dynamic programming* to estimate the utilities of certain actions in certain states**

both approaches can be- and often are combined,
but we henceforth solely focus on the second one

further fancy models of intelligent agents



remarks

yes! there is a philosophical ambiguity . . .

some people think of states as describing an agent's environment

other people think of states as describing an agent's internal condition

from a practical point of view, this does not matter . . .

as long as we are aware of what we consider as a state and then stick to our modeling choice (when coding), we can use the same upcoming math to reason about and implement decision making agents

⇔ **we can and henceforth will simply talk about “the state an agent is in”**

 **note**

both above diagrams describe agents who interact with an environment by executing

actions a

which may change

states s

and also trigger

rewards r

but what about observations o ?

let's simplify our lives by making
several simplifying assumptions ...

simplifying assumptions (1)

we henceforth assume that

observations o provide complete information about states s

\Leftrightarrow observations o *determine* states \Leftrightarrow states s are *observable*

\Rightarrow **we can ignore observations and simply reason about states**

observability is a term from [control theory](#)

states not fully described by observations are called *partially observable* states but, for now, we won't worry about those . . .

simplifying assumptions (2)

we henceforth assume that

the sets of states and actions

$$\mathcal{S} = \{s_1, s_2, \dots\}$$

$$\mathcal{A} = \{a_1, a_2, \dots\}$$

are both discrete

rewards are real numbers $r \in \mathcal{R} \subset \mathbb{R}$

the sets \mathcal{S} , \mathcal{A} , and \mathcal{R} are of finite sizes

(much) later, we will relax these assumptions

 **note**

behavior is best seen as a *stochastic* temporal phenomenon

⇒ we should think of states, actions, and rewards as *random variables*

⇔ at each time step t , the decision making agent finds itself in a state

$$S_t \in \mathcal{S}$$

selects and performs an action

$$A_t \in \mathcal{A}$$

upon which it receives a reward

$$R_{t+1} \in \mathcal{R}$$

and transits to

$$S_{t+1} \in \mathcal{S}$$

yet, when we run a decision making agent, the rvs S_t , A_t , and R_t will “collapse” to certain particular values

to refer to such *instantiations* of the rvs S_t , A_t , and R_t , we write

$s[t] \in \mathcal{S} \quad \Leftrightarrow \quad$ particular state an agent is in at time t

$a[t] \in \mathcal{A} \quad \Leftrightarrow \quad$ particular action an agent takes at time t

$r[t + 1] \in \mathcal{R} \quad \Leftrightarrow \quad$ particular reward an agent receives at time $t + 1$

 **note**

our simplifying assumptions allow us to study the reinforcement learning problem from the perspective of (finite) *Markov decision processes (MDPs)*

MDPs provide a mathematical formalization of sequential decision making

they provide an idealized mathematical model of the reinforcement learning problem which allows for precise theoretical statements

as usual in AI, there is a gap between the model's mathematical tractability and the breadth of its (real world) applicability

the last statement basically says that phenomena of intelligence are too complex for mathematical modeling; w.r.t. reinforcement learning it means that an MDP is but a **useful *approximative model*** of what we really would want to describe . . .

Markov decision process (MDP)

\Leftrightarrow a tuple $(\mathcal{S}, \mathcal{A}, T, R)$ where

$\mathcal{S} \equiv$ set of states

$\mathcal{A} \equiv$ set of actions

$T(s', s, a) \equiv$ state transition probabilities

$$= p(S_{t+1} = s' \mid S_t = s, A_t = a)$$

$R(s, a) \equiv$ expected immediate rewards / reinforcements

$$= \mathbb{E}[R_{t+1} \mid S_t = s, A_t = a]$$



here is how we can make sense of the notion of *expected* immediate rewards

$$R(s, a) = \mathbb{E}[R_{t+1} \mid S_t = s, A_t = a] = \sum_{r \in \mathcal{R}} r \cdot p(r \mid s, a)$$

but what is r ? and what is $p(r \mid s, a)$? neither occurs in the MDP definition !!!

so, let's consider this ...

 **note**

if an agent could compute a function $Util(s)$ to assess a state's *utility*, then we could let

$$r = Util(s')$$

and think of the expected immediate reward w.r.t. s and a as follows

$$R(s, a) = \mathbb{E}[R_{t+1} \mid S_t = s, A_t = a] = \sum_{s' \in \mathcal{S}} Util(s') \cdot p(s' \mid s, a)$$

such a function $Util(s)$ is not included in most common MDP definitions

it may seem crazy to say this, but “theoretically” $Util(s)$ is not necessary

we brought it up, as it ties back to our earlier discussions and (hopefully) removes a point of great confusion for RL novices ...



let's dwell on the point that $Util(s)$ is “theoretically” unnecessary

one can develop the (MDP-based) theory of RL and most RL algorithms that we know of without having to worry about how to implement $R(s, a)$

indeed, this is what most textbooks, lectures, blogs, vlogs we know of do

however, as soon as it comes to practice and coding, the question of how to implement / represent / realize $R(s, a)$ becomes critical

in fact, so called *reward shaping* is a crucial success factor in applied RL

next, we proceed as in most textbooks but at least we now already know about the need for *reward shaping*

 **note**

if $|\mathcal{S}|$ and $|\mathcal{A}|$ are finite, the transition function T and reward function R can be implemented as (possibly humongous) arrays or *lookup tables*

\Leftrightarrow if $|\mathcal{S}|$ and $|\mathcal{A}|$ are finite, then both functions T and R can be *tabulated*

in our following discussion, we will pretend these tables are available, i.e. that somebody has provided us with all the numbers they contain

for our initial theoretical developments and considerations, this is OK in most practical settings, however, it is not

(much) later on, we will see how to handle such real world settings ...

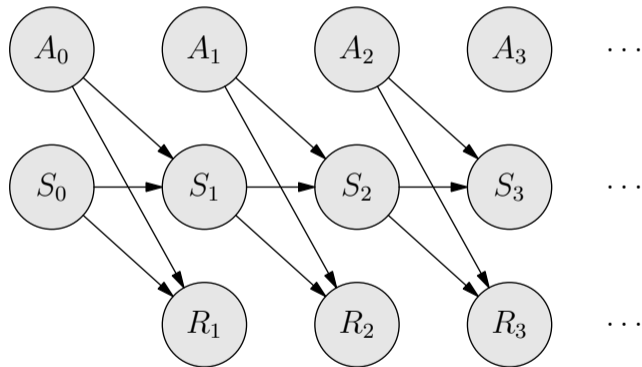
 **note**

alternatively, we could have followed the definition in *the RL textbook* by **Sutton** and **Barto** and said

an MDP is a stochastic process over rvs S_t , A_t , and R_t where

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

Bayesian network representation



this Bayes net definition subsumes the one we presented before, since we have

$$T(s', s, a) = p(S_{t+1} = s' \mid S_t = s, A_t = a) = \sum_{r \in \mathcal{R}} p(s', r \mid s, a)$$

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

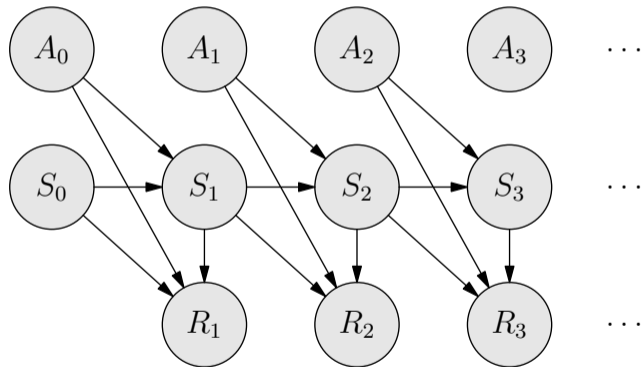
this explicitly resolves confusions as to expected rewards ...

some authors also think of the reward function as a three-argument function

$$R'''(s', s, a) = \mathbb{E}[R_{t+1} \mid S_{t+1} = s', S_t = s, A_t = a] = \sum_{r \in \mathcal{R}} r \cdot \frac{p(s', r \mid s, a)}{p(s' \mid s, a)}$$

this is the most general view on MDPs allowing for the most flexible modeling
however, —to our knowledge— it is rarely necessary in practical applications

Bayesian network representation



some authors also think of the reward function as a univariate function $R'(s)$

this is not unreasonable because we can relate the idea to our own definition

$$\begin{aligned} R(s, a) &= \mathbb{E}[R_{t+1} \mid S_t = s, A_t = a] \\ &= \sum_{r \in \mathcal{R}} r \cdot p(r \mid s, a) && (r \text{ and } p(r|s, a) \text{ are again but auxiliary}) \\ &= \sum_{r \in \mathcal{R}} \sum_{s' \in \mathcal{S}} r \cdot p(r \mid s, a, s') \cdot p(s' \mid s, a) \\ &= \sum_{s' \in \mathcal{S}} \sum_{r \in \mathcal{R}} r \cdot p(r \mid s, a, s') \cdot p(s' \mid s, a) \\ &\equiv \sum_{s' \in \mathcal{S}} R'(s') \cdot p(s' \mid s, a) \end{aligned}$$

if you have not yet seen marginalization as on the previous slide, then ...

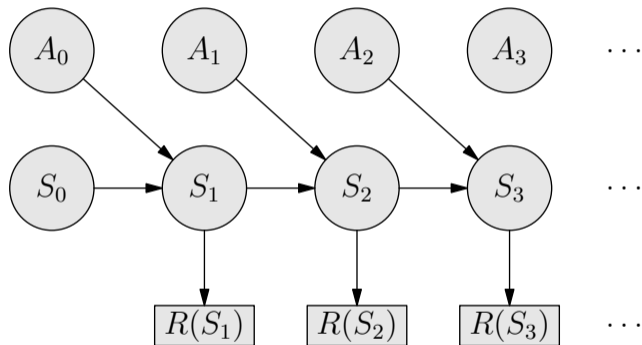
observe that

$$\begin{aligned} p(X | Z) &= \frac{p(X, Z)}{p(Z)} = \sum_y \frac{p(X, Y = y, Z)}{p(Z)} = \sum_y \frac{p(X | Y = y, Z) p(Y = y, Z)}{p(Z)} = \sum_y \frac{p(X | Y = y, Z) p(Y = y | Z) p(Z)}{p(Z)} \\ &= \sum_y p(X | Z, Y = y) p(Y = y | Z) \end{aligned}$$

equivalently

$$\begin{aligned} p(X | Z) &= \sum_y p(X, Y = y | Z) = \sum_y \frac{p(X, Y = y, Z)}{p(Z)} = \sum_y \frac{p(X | Y = y, Z) p(Y = y, Z)}{p(Z)} = \sum_y \frac{p(X | Y = y, Z) p(Y = y | Z) p(Z)}{p(Z)} \\ &= \sum_y p(X | Z, Y = y) p(Y = y | Z) \end{aligned}$$

Bayesian network representation



circles denote rvs, boxes denote deterministic functions of rvs

remarks

“as usual”, the Bayesian network view on MDPs is arguably more powerful

however, the majority of the literature (still) follows the state transition view

in fact, here is a quote from the 2nd edition of the **Sutton** and **Barto** book

Our characterization of the dynamics of an MDP in terms of $p(s', r | s, a)$ is slightly unusual. It is more common in the MDP literature to describe the dynamics in terms of the state transition probabilities $p(s' | s, a)$ and expected next rewards $r(s, a)$. In reinforcement learning, however, we more often have to refer to individual actual or sample rewards (rather than just their expected values). Our notation also makes it plainer that S_t and R_t are in general jointly determined, and thus must have the same time index. In teaching reinforcement learning, we have found our notation to be more straightforward conceptually and easier to understand.

remarks

it is what it is . . .

reinforcement learning is a comparatively new discipline and its practitioners still struggle for a standardized notation and a common point of view

this poses one (of the many) challenges for people who are just beginning to learn about RL . . .

further terminology

if we *sample* an MDP, we obtain an (infinite) **trajectory**

$$s[0], a[0], r[1], s[1], a[1], r[2], s[2], a[2], r[3], s[3], a[3], \dots$$

a trajectory of finite length is called an **episode** or a **trial**

$$s[0], a[0], r[1], s[1], a[1], \dots, r[h], s[h], a[h]$$

a tuple $(s[t], a[t], r[t + 1], s[t + 1])$ is said to be a **transition**

finally, a sequence of transitions is also called a **roll-out**

example episodes (all starting in the same state)



[3, 1], →



+10, [4, 1], ?



[3, 1], ←



+5, [2, 0], ←



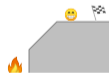
-10, [1, 0], ?



[3, 1], ←



+5, [2, 0], →



-1, [3, 0], ←



-1, [2, 0], →



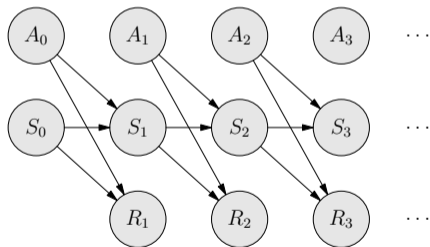
-1, [3, 0], →



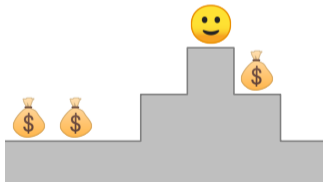
+10, [4, 0], ?



R_{t+1} is explicitly computed from S_t and A_t but it implicitly also depends on S_{t-1} and A_{t-1} because these are required to determine $S_t \dots$



- ⇒ MDPs can be seen as a formalization of sequential decision making in which actions (state transitions) not only impact immediate rewards but also future states and thus future rewards
- ⇔ MDPs give rise to the notion of *delayed rewards* and the need to trade off immediate- and delayed rewards



given this, we can now state the ...

goal of a reinforcement learning agent

⇔ maximize *total amount of reward*

⇔ maximize not immediate reward but long term *cumulative reward*

note: the use of *reward signals* to formalize learning objectives distinguishes reinforcement learning from other machine learning paradigms

question

how to model cumulative rewards ?

answer

most commonly, it is done like this . . .

for learning **episodic tasks** which start in a state $S_t = s$ and end after / have a *horizon* of $h < \infty$ steps, we may simply consider the (expected) future **return**

$$G_t = R_{t+1} + R_{t+2} + R_{t+3} + \dots + R_{t+h}$$

learning to play a game poses an episodic task

e.g. no game of *tic tac toe* lasts longer than 9 moves

i.e. every game has final states where episodes end

for learning **continuing tasks** which last forever ($h \rightarrow \infty$), the above will likely not work (as returns might grow beyond all bounds and could thus not be maximized)

the common remedy is to consider the (expected) future **discounted return**

$$G_t = R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots = \sum_{\tau=0}^{\infty} \gamma^{\tau} R_{t+1+\tau}$$

where the parameter $0 \leq \gamma < 1$ is called the **discount rate** or **discount factor**

control tasks are examples of continuing tasks

 **note**

there is a hidden recursion because

$$\begin{aligned}G_t &= R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \gamma^3 R_{t+4} + \dots \\&= R_{t+1} + \gamma [R_{t+2} + \gamma R_{t+3} + \gamma^2 R_{t+4} + \dots] \\&= R_{t+1} + \gamma G_{t+1}\end{aligned}$$

so that (with $\gamma \neq 0$)

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

in the literature, we often read

the discount factor $0 \leq \gamma < 1$ is used to assign more importance to immediate rewards than to later ones

indeed, a reward received τ steps in the future is only worth a fraction of $\gamma^{\tau-1}$ of what it would be worth in the present

moreover, if $\gamma = 0$ (so that $0^0 = 1$ by definition and $0^\tau = 0$ for all $\tau > 0$), an RL agent will only maximize immediate returns

an agent that learns with $\gamma = 0$ is called **nearsighted / myopic**
an agent that learns with $\gamma \approx 1$ is called **farsighted / hyperopic**

but let's be honest . . .

 **note**

discounting is but a “trick” to bound the above infinite sum

because, if $|R_{t+1}| \leq r_{\max} \in \mathbb{R}$ at all times t and if $0 \leq \gamma < 1$, then

$$\begin{aligned} \sum_{\tau=0}^{\infty} \gamma^{\tau} R_{t+1+\tau} &\leq \sum_{\tau=0}^{\infty} \gamma^{\tau} r_{\max} \\ &= r_{\max} \sum_{\tau=0}^{\infty} \gamma^{\tau} \\ &= r_{\max} \frac{1}{1-\gamma} \end{aligned}$$

recall (our very important) exercise 1.1

the series of powers of the discount factor γ

$$\sum_{\tau=0}^{\infty} \gamma^{\tau}$$

is but the specific geometric series with $a = 1$

\Rightarrow we do indeed have the rather pivotal equality

$$\sum_{\tau=0}^{\infty} \gamma^{\tau} = \frac{1}{1-\gamma}$$

 **note**

observe that (expected) future discounted returns are all about the future and do not care about the past

to simplify notation, we henceforth assume that $t = 0$ denotes the present in which the agent is in a known state $S_0 = s$ so that we may simply write

$$G_0 = \sum_{t=0}^{\infty} \gamma^t R_{t+1} = \sum_{t=0}^{\infty} \gamma^t R(S_t, A_t)$$

this expression is also called **infinite-horizon discounted model** with $0 \leq \gamma < 1$

(it subsumes finite-horizon models ($h < \infty$) w/o discount if we allow $\gamma = 1$ and $\forall t > h : R_{t+1} = 0$)

question

but what about maximizing returns ?

question

but what about maximizing returns ?

answer

we haven't talked about *policies* yet . . .

we therefore recall from lecture 01 . . .

a **stochastic policy** is a *probabilistic map*

$$\pi : \mathcal{S} \times \mathcal{A} \rightarrow [0, 1]$$

with $0 \leq \pi(s, a) \equiv \pi(a | s) \equiv p(a | s) \leq 1$ and $\sum_a \pi(a | s) = 1$

 **note**

⇒ a stochastic policy

$$\pi(a | s) = p(A_t = a | S_t = s)$$

is a distribution over actions given states

MDP policies only regard / involve the current state, not a history of states

throughout, we will assume that policies are *stationary* (time-independent)

$$\forall t : A_t \sim \pi$$

question

why policies ?

answer

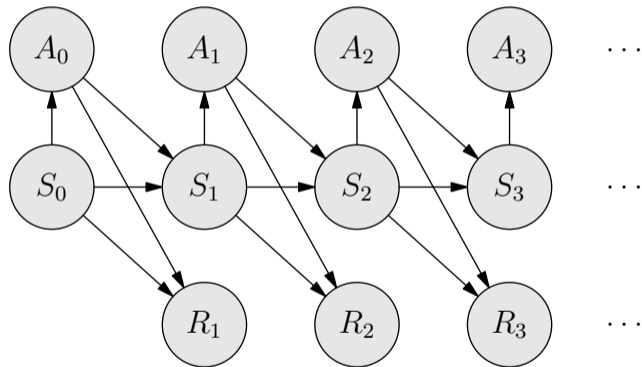
if an agent follows a policy, it has a criterion for deciding “*what* to do *when*”

actions are not selected arbitrarily but in “some” informed manner based on “some” background knowledge / experience

informed decision making provides “some” control over long-term returns

introducing *stochastic policies*, our MDP model becomes . . .

Bayesian network representation



 **note**

the problem at the heart of the MDP framework of reinforcement learning is to *learn* a policy which specifies which action a an agent should choose in state s

but before we can study *policy learning*, we must consider how (stochastic) policies impact our previous considerations . . .

above, we assumed we know $S_t = s$ and $A_t = a$

knowing s and a allows us to retrieve the value

$$R(s, a) = \mathbb{E}[R_{t+1} \mid S_t = s, A_t = a]$$

from a table

but what if we only know $S_t = s$ and $A_t \sim \pi(a \mid s)$?

in other words, what is “ $\mathbb{E}[R_{t+1} \mid S_t = s, A_t \sim \pi]$ ” ?

if all we know is $A_t \sim \pi(a | s)$, we can average it out to get

$$\begin{aligned}\mathbb{E}_{A_t \sim \pi} \left[\mathbb{E}[R_{t+1} | S_t = s, A_t] \right] &\equiv \mathbb{E}_{\pi} \left[\mathbb{E}[R_{t+1} | S_t = s, A_t] \right] \\ &= \mathbb{E}_{\pi} \left[\sum_r r \cdot p(r | s, A_t) \right] \\ &= \sum_a \left[\sum_r r \cdot p(r | s, a) \right] \cdot \pi(a | s) \\ &= \sum_a \sum_r r \cdot \frac{p(r, s, a)}{p(s, a)} \cdot \frac{p(s, a)}{p(s)} \\ &= \sum_a \sum_r r \cdot p(r, a | s) \\ &= \sum_r r \cdot p(r | s) \\ &= \mathbb{E}[R_{t+1} | S_t = s]\end{aligned}$$

to retrace what we just did, recall the **law of total expectation** (or *tower law*)

$$\begin{aligned}\mathbb{E}[X] &= \sum_x x \cdot p(X = x) \\ &= \sum_x x \cdot \sum_y p(X = x, Y = y) \\ &= \sum_x x \cdot \sum_y p(X = x | Y = y) \cdot p(Y = y) \\ &= \sum_y p(Y = y) \cdot \sum_x x \cdot p(X = x | Y = y) \\ &= \sum_y p(Y = y) \cdot \mathbb{E}[X | Y = y] \\ &= \mathbb{E}[\mathbb{E}[X | Y]]\end{aligned}$$

while we are at it . . .

if X is conditionally independent (*ci*) of Y_t given Y_{t+1}

$$X \perp\!\!\!\perp Y_t \mid Y_{t+1} \Leftrightarrow p(X \mid Y_{t+1}, Y_t) = p(X \mid Y_{t+1})$$

then we also have the (soon to be important) expectation

$$\begin{aligned} \mathbb{E}[X \mid Y_t = y] &= \sum_x x \cdot p(X = x \mid Y_t = y) \\ &= \sum_x x \cdot \sum_{y'} p(X = x \mid Y_{t+1} = y', Y_t = y) \cdot p(Y_{t+1} = y' \mid Y_t = y) \\ &\stackrel{ci}{=} \sum_x x \cdot \sum_{y'} p(X = x \mid Y_{t+1} = y') \cdot p(Y_{t+1} = y' \mid Y_t = y) \\ &= \sum_{y'} \sum_x x \cdot p(X = x \mid Y_{t+1} = y') \cdot p(Y_{t+1} = y' \mid Y_t = y) \\ &= \sum_{y'} \mathbb{E}[X \mid Y_{t+1} = y'] \cdot p(Y_{t+1} = y' \mid Y_t = y) \end{aligned} \tag{*}$$

 **note**

to denote an expectation $\mathbb{E}_{A_t \sim \pi} [\cdot]$ w.r.t. a policy $\pi(a | s)$,
we already wrote and henceforth always will simply write

$$\mathbb{E}_{A_t \sim \pi} [\cdot] = \mathbb{E}_{\pi} [\cdot]$$

almost everybody does this

question

if stochastic policies impact immediate rewards, then what about long-term returns ?

answer

we need to extend the notion of infinite horizon discounted model to *value functions*

value function / state-value function

- ⇔ a function $V_\pi : \mathcal{S} \rightarrow \mathbb{R}$ that computes the *value of a state* under a given policy π
- ⇔ a function $V_\pi : \mathcal{S} \rightarrow \mathbb{R}$ that computes the expected return when starting in $s \in \mathcal{S}$ and following policy π thereafter

$$V_\pi(s) = \mathbb{E}_\pi \left[\sum_{t=0}^{\infty} \gamma^t R_{t+1} \mid S_0 = s \right] = \mathbb{E}_\pi \left[\sum_{t=0}^{\infty} \gamma^t R(S_t, A_t) \mid S_0 = s \right]$$

quality function / state-action-value function

\Leftrightarrow a function $Q_\pi : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ that computes the expected return starting from s , taking action a , and then following policy π

$$Q_\pi(s, a) = \mathbb{E}_\pi \left[\sum_{t=0}^{\infty} \gamma^t R_{t+1} \mid S_0 = s, A_0 = a \right]$$



V_π and Q_π are interrelated

$$V_\pi(s) = \sum_a Q_\pi(s, a) \pi(a | s)$$

moreover ...

previously, we wrote

$$G_0 = \sum_{t=0}^{\infty} \gamma^t R_{t+1}$$

and have seen that

$$G_0 = R_1 + \gamma G_1$$

⇒ we have

$$\begin{aligned}V_{\pi}(s) &= \mathbb{E}_{\pi}[G_0 \mid S_0 = s] \\&= \mathbb{E}_{\pi}[R_1 + \gamma G_1 \mid S_0 = s] \\&= \mathbb{E}_{\pi}[R_1 \mid S_0 = s] + \mathbb{E}_{\pi}[\gamma G_1 \mid S_0 = s] \\&= \mathbb{E}_{\pi}[R_1 \mid S_0 = s] + \gamma \mathbb{E}_{\pi}[G_1 \mid S_0 = s]\end{aligned}$$

the second term resembles a value function but is conditioned on S_0 instead on S_1

however ...

using (★) on slide 55, we find

$$\begin{aligned}\mathbb{E}_\pi[G_1 \mid S_0 = s] &= \sum_{s'} \mathbb{E}_\pi[G_1 \mid S_1 = s'] \cdot p(S_1 = s' \mid S_0 = s) \\ &= \sum_a \sum_{s'} \mathbb{E}_\pi[G_1 \mid S_1 = s'] \cdot p(s' \mid s, a) \cdot \pi(a \mid s) \\ &= \sum_a \sum_{s'} V_\pi(s') \cdot p(s' \mid s, a) \cdot \pi(a \mid s)\end{aligned}$$

most bloggers / vloggers don't understand this step and get it wrong

⇒ since

$$\mathbb{E}_\pi[R_1 \mid S_0 = s] = \sum_a R(s, a) \pi(a \mid s)$$

and

$$\mathbb{E}_\pi[G_1 \mid S_0 = s] = \sum_a \sum_{s'} V_\pi(s') p(s' \mid s, a) \pi(a \mid s)$$

we find

$$\begin{aligned} V_\pi(s) &= \mathbb{E}_\pi[G_0 \mid S_0 = s] \\ &= \mathbb{E}_\pi[R_1 \mid S_0 = s] + \gamma \mathbb{E}_\pi[G_1 \mid S_0 = s] \\ &= \sum_a \left[R(s, a) + \gamma \sum_{s'} V_\pi(s') p(s' \mid s, a) \right] \cdot \pi(a \mid s) \end{aligned}$$

Bellman equation for $V_\pi(s)$

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

similarly, we have

$$Q_{\pi}(s, a) = \mathbb{E}_{\pi}[G_0 \mid S_0 = s, A_0 = a] = \mathbb{E}_{\pi}[R_1 + \gamma G_1 \mid S_0 = s, A_0 = a] = \underbrace{\mathbb{E}_{\pi}[R_1 \mid S_0 = s, A_0 = a]}_{=R(s,a)} + \gamma \mathbb{E}_{\pi}[G_1 \mid S_0 = s, A_0 = a]$$

and we note that

$$\begin{aligned} \mathbb{E}_{\pi}[G_1 \mid s, a] &= \sum_{s'} \mathbb{E}_{\pi}[G_1 \mid s', s, a] p(s' \mid s, a) = \sum_{s'} \sum_{a'} \mathbb{E}_{\pi}[G_1 \mid s', s, a, a'] p(s' \mid s, a) p(a' \mid s', s, a) \\ &\stackrel{ci}{=} \sum_{s'} \sum_{a'} \mathbb{E}_{\pi}[G_1 \mid s', a'] p(s' \mid s, a) p(a' \mid s') \\ &= \sum_{s'} \sum_{a'} Q_{\pi}(s', a') p(s' \mid s, a) p(a' \mid s') \end{aligned}$$

we therefore have

$$Q_{\pi}(s, a) = R(s, a) + \gamma \sum_{s'} \sum_{a'} Q_{\pi}(s', a') p(s' \mid s, a) p(a' \mid s')$$

Bellman equation for $Q_\pi(s, a)$

$$Q_\pi(s, a) = R(s, a) + \gamma \sum_{s'} \sum_{a'} Q_\pi(s', a') \pi(a' | s') p(s' | s, a)$$

note: this differs from what you will find on most websites
however, it agrees with *our* modeling assumptions !

 **note**

the *real* problem at the heart of MDPs is to determine an **optimal policy**

$$\pi_* : \mathcal{S} \times \mathcal{A} \rightarrow [0, 1]$$

such that action selection according to $\pi_*(a \mid s)$ entails optimal behavior

this statement is not yet precise,
so that we will get back to it later
for now, we point out . . .

a policy π is **better** than a policy π' , written as $\pi \succcurlyeq \pi'$, if its expected return is always *greater than or equal to* that of π'

$$\pi \succcurlyeq \pi' \quad \Leftrightarrow \quad \forall s : V_{\pi}(s) \succcurlyeq V_{\pi'}(s)$$

 **note**

this induces a partial ordering over the set of all policies (if one policy is better than another at one state but worse elsewhere, the two policies are incomparable)

note / recall that every partially ordered set has *at least* one maximal element

⇒ there must exist at least one policy better than all others, $\exists \pi_* : \forall s : V_{\pi_*}(s) \geq V_{\pi}(s)$

⇒ **for any Markov decision process, there exists at least one optimal policy π_***

in case there are several optimal policies, it is traditional to simply call them all π_*

optimal value function

⇔ the maximum value function over all policies

$$V_*(s) = \max_{\pi} V_{\pi}(s)$$

optimal quality function

⇔ the maximum quality function over all policies

$$Q_*(s, a) = \max_{\pi} Q_{\pi}(s, a)$$

optimal policy

given an optimal quality function, we can compute

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

Theorem (w/o proof)

for any MDP

there exists an optimal policy better than all other policies

$$\exists \pi_* : \forall \pi : \pi_* \geq \pi$$

all optimal policies achieve the optimal value function

$$V_{\pi_*} = V_*$$

all optimal policies achieve the optimal quality function

$$Q_{\pi_*} = Q_*$$

Theorem

an MDP can have several optimal policies π_ , π'_* , ... but only one optimal value function V_**

Proof (by contradiction).

suppose there are two optimal policies π_* and π'_* which have two *distinct* optimal value functions V_* and V'_*

if V_* and V'_* are distinct, there must exist a state s with $V_*(s) \neq V'_*(s)$ and we may assume w.l.o.g. that $V_*(s) < V'_*(s)$

however, if $V_*(s)$ is less than $V'_*(s)$, then π_* can't be optimal as it would be better to follow π'_*

\Rightarrow if π_* and π'_* are both optimal, V_* and V'_* must be the same

□

note: we will see examples in the exercises !!!

 **note**

there are two fundamental questions when working with MDPs

- 1) **how to *find an optimal policy* given a model $(\mathcal{S}, \mathcal{A}, T, R)$?**
- 2) **how to *learn an optimal policy* without having a model ?**

in the next couple of lectures, we address both questions

summary

we now know about

the reinforcement learning problem

Markov decision processes (MDPs)

value functions and policies

Bellman equations